

1968

Dispersion effects in industrial property life analysis

George Emmett Lamp Jr.
Iowa State University

Follow this and additional works at: <https://lib.dr.iastate.edu/rtd>



Part of the [Industrial Engineering Commons](#)

Recommended Citation

Lamp, George Emmett Jr., "Dispersion effects in industrial property life analysis " (1968). *Retrospective Theses and Dissertations*. 3251.
<https://lib.dr.iastate.edu/rtd/3251>

This Dissertation is brought to you for free and open access by the Iowa State University Capstones, Theses and Dissertations at Iowa State University Digital Repository. It has been accepted for inclusion in Retrospective Theses and Dissertations by an authorized administrator of Iowa State University Digital Repository. For more information, please contact digirep@iastate.edu.

This dissertation has been
microfilmed exactly as received

68-10,471

LAMP, Jr., George Emmett, 1933-
DISPERSION EFFECTS IN INDUSTRIAL
PROPERTY LIFE ANALYSIS.

Iowa State University, Ph. D., 1968
Engineering, industrial

University Microfilms, Inc., Ann Arbor, Michigan

© George Emmett Lamp, Jr. 1968

DISPERSION EFFECTS IN INDUSTRIAL PROPERTY LIFE ANALYSIS

by

George Emmett Lamp, Jr.

A Dissertation Submitted to the
Graduate Faculty in Partial Fulfillment of
The Requirements for the Degree of
DOCTOR OF PHILOSOPHY

Major Subject: Engineering Valuation

Approved:,,

Signature was redacted for privacy.

In Charge of Major Work

Signature was redacted for privacy.

~~Head of Major Department~~

Signature was redacted for privacy.

Dean of Graduate College

Iowa State University
Of Science and Technology
Ames, Iowa

1968

TABLE OF CONTENTS

	Page
INTRODUCTION	1
Depreciation	1
Life Estimation	4
Life Analysis	5
OBJECTIVES OF INVESTIGATION	20
PRESENT ACTUARIAL METHODS OF LIFE ANALYSIS	22
Related Concepts	23
Selection and Aggregation of Property Data	37
Original Life Table	40
Individual-unit method	42
Original-group method	42
Composite original-group method	43
Multiple original-group method	44
Annual-rate method	45
Methods of Obtaining a Smoothed Life Table	46
Judgment method	48
Matching method	48
Statistical curve fitting methods	51
INVESTIGATION	59
Simulation of Retirement Ratios	59
Normal Approximation	61
RESULTS OF THE INVESTIGATION	67
A PROCEDURE FOR FITTING A POLYNOMIAL TO RETIREMENT RATIOS	88
Assumptions	88
Development of Procedure	90
The Procedure	94
Comments	95

	Page
DISCUSSION	100
Estimators of the Composite Retirement Ratios	104
Estimators of the Polynomial Coefficients	106
Effect of Dollars on Computing Variances	111
EXAMPLE	113
CONCLUSIONS	131
LITERATURE CITED	133
ACKNOWLEDGMENTS	136
APPENDIX A - GENERAL FLOW CHART OF SIMULATION PROGRAM	137
APPENDIX B - GENERAL FLOW CHART OF NORMAL APPROXIMATION PROGRAM	140
APPENDIX C - MAXIMUM-LIKELIHOOD ESTIMATORS OF THE POLYNOMIAL COEFFICIENTS	143
APPENDIX D - PRELIMINARY APPROACH	150
APPENDIX E - GENERAL FLOW CHART OF PROGRAM TO IMPLEMENT THE PROCEDURE	156
APPENDIX F - TESTING FOR NORMALITY	161

INTRODUCTION

Businesses may utilize estimates of the mortality behavior of property for a number of purposes, such as computing income tax liability, computing the rate base and depreciation expenses for rate regulation, and making management decisions relating to property. A usage of estimates of mortality behavior of property common to all of these purposes is in the calculation of depreciation. Depreciation calculations generally require an estimate of the probable average service life of the property group, or the probable service life of the unit of property, and may require an estimate of the probable retirement dispersion pattern of the property group. The process of estimating the probable average service life or probable service life and, if feasible, the probable retirement dispersion pattern is called life estimation. An extensive knowledge of the past mortality behavior of the same or a similar property forms a useful part of the information for life estimation. The process of aggregating and analyzing historical data to obtain this knowledge is called life analysis.

Depreciation

W. C. Fitch, after an extensive study, formulated a general definition of depreciation (6, p. 76):

Depreciation is the decrease in the number of available units of service which a unit of property or group of property units can be expected to render.

Three basic concepts of depreciation are frequently recognized: cost, value, and physical condition (6, p. 10; 20, p. 175). Fitch states formal

definitions of cost-depreciation and value-depreciation (6, pp. 76-77):

Cost-depreciation is the decrease in the available units of service expressed as a function of the cost of the property.

Value-depreciation is the change in the present worth of the anticipated returns from the services to be rendered by a property.

and summarizes all three basic concepts as (6, p. 10):

Cost-depreciation is the allocation of the purchase price over the life of the equipment. Value-depreciation is the change in anticipated benefits between two points in time. . . . physical condition is an estimate of the percent of the tangible decay of a property.

Bonbright mentions a fourth concept of depreciation (2, p. 185):

" . . . the difference between the present worth of the old and obsolescent asset and the present worth of the hypothetical, new and modern asset."

The Federal and State governments generally require the use of the cost concept of depreciation for the purpose of estimating income tax liability. Rate regulation agencies may utilize the cost or physical condition or difference in value concepts and/or some combination(s) of the four concepts of depreciation (10, p. 29). The pertinent constitution, statutes, and court decisions and the policies and decisions of the particular regulatory agency may prescribe which concept(s) is appropriate for the agency and the regulated business to use for regulatory purposes. Management may use whichever concept or combination of concepts they deem appropriate for making a particular management decision.

Depreciation calculations generally require, to a greater or lesser extent, estimates of one or more of the mortality characteristics of the property. The mortality characteristics specifically referred to are the probable average service life of a property group, or the probable service

life of a unit of property, and the retirement or mortality dispersion pattern of a property group. Winfrey defined probable service life and probable average service life as (28, p. 12):

The probable service life of an individual unit is that period of time extending from its date of installation to the forecasted date when it probable will be retired.

The probable average service life of a group of individual units is the average of the probable service lives of the units of the group.

The retirement dispersion pattern refers to the distribution of the ages at retirement of the units comprising the property group. Probable average service life can be calculated from the probable retirement dispersion pattern; the reverse is not true.

Cost-depreciation requires estimates of one or more of the mortality characteristics of the property. Value-depreciation does not directly require estimates of any of the mortality characteristics of the property; however, the process of estimating the anticipated benefits may utilize estimates of one or more of the mortality characteristics. Physical condition is estimated, generally, by an inspection of the property (20, p. 178). Therefore, estimates of the mortality characteristics are not directly involved in estimating depreciation in the sense of physical condition.

The word depreciation will be used in the sense of cost-depreciation in the remainder of this dissertation, unless otherwise noted, to simplify the discussion.

The units of service which a property can be expected to render are generally measured in terms of years of service or units of production.

Years of service are the most frequently used measure (4, p. 30) and are used as the measure of service life in this dissertation.

Annual depreciation is that portion of the cost of a property charged as an expense (and, hence, charged against revenue) for a year. Accrued depreciation as of a given date is the total depreciation of the unretired property charged as an expense from the time of installation of the property until that date.

A more extensive treatment of depreciation may be found in Fitch (6), Grant and Norton (8), and Marston, et al. (20).

Life Estimation

The process of life estimation can be divided into two parts. The first part is the collection of relevant information. The second part is the application of expert judgment to the information available to estimate the mortality behavior of the property.

Relevant information includes, but is not necessarily limited to, the results of a life analysis and analyses of economic trends, technological progress, and policies and decisions of governmental bodies and agencies and of management. While trends based on historical information can usually be extended and extrapolated, a degree of uncertainty is present in any attempt to predict the future; hence, expert judgment is an essential part of life estimation.

While the author strongly recommends the development and use of retirement data and survivor curves as the basis of estimating the probable life of property units, he does not mean to infer that expert judgment should be done away with in favor of pure statistical treatment. Each individual item, each group of items, and each property or company must be dealt with in the light of its present condition, its character and amount of service or production, and

its relation to the present and probable future economic trends, art of manufacture, and management policies. Tables of probable service lives, type survivor curves, and statistical methods are simply means of recording past experience to use in predicting what the future service might be (28, p. 9).

Life Analysis

Life analysis is the process of aggregating and analyzing the historical record of property for the purpose of obtaining information about the mortality characteristics of the property. Life analysis and life estimation are different processes since the former is concerned with an analysis of the past whereas the latter is generally concerned with a prediction of the future. The end result of a life analysis is an estimate of the probable average service life and, if possible, of the probable retirement dispersion pattern, as well as a knowledge of any discernable trends in either, experienced or being experienced by the property under study.

The plant property records are a primary source of data for studying the past mortality experience of property. A separate account may be kept for each individual unit of property or two or more individual units may be combined into a group and a record kept of the units as a group. A group account in which the installations in a single year of a given type(s) of property are recorded is called a vintage account and the group of units is called a vintage group. A group account in which the installations (of the same type or types of property) of successive years are recorded is called a continuous or "open-end" account. The ensuing discussion is based on the life analysis of group property.

A complete property record would permit determination of at least the following:

1. The amount of property installed each year (i.e., the amount installed each year as a vintage group),
2. The age at retirement of the property already retired from each vintage group, and
3. The total amount of property in each vintage group surviving at the beginning of each year (plant balance of each vintage group at the beginning of each year).

A particular property record may not contain all of the above information. Sometimes the only information available is the amount of property installed each year, the amount retired each year, and the total plant balance each year.

The extent of the property data available affects the choice of methods of analyzing the data. The statistical methods of life analysis are often divided into two categories: the turnover methods and the actuarial methods. The turnover methods require data on the amount of property installed each year, the amount retired each year, and the total plant balance each year. The actuarial methods generally require a complete property record (as described in the preceding paragraph).

The turnover methods, with one exception, yield only an indication of the probable average service life. The simulated plant balance method, often classified as a turnover method since the data requirements are similar, does yield estimates of both probable average service life and probable retirement dispersion pattern.

Estimates of both the probable average service life and the probable retirement dispersion pattern can be obtained by use of the actuarial methods. The tabulation of the raw data frequently results in an incomplete, original life table. A life table is the amount, percent or proportion of property surviving at each age; an original life table is a life table calculated from the observed data. Before the probable average service life and the probable retirement dispersion pattern can be estimated, the original life table generally must be smoothed and extended to zero survivors or zero percent surviving. Even if the original life table is complete, the common practice is to smooth the original life table and use the interpolated values to estimate the probable average service life and probable retirement dispersion pattern.

A commonly used technique of smoothing and of extending (if necessary) the life table is by fitting a multiple linear regression equation, generally a polynomial, to the retirement ratios by the method of least-squares (4, p. 5). The retirement ratio for an age interval is the amount of property, in terms of proportion, percent, units, or dollars, retired during the age interval divided by the amount of property surviving at the beginning of the age interval. A smoothed life table can be calculated from the retirement ratio polynomial by starting with the amount installed (1.00, 100%, units, or dollars) and successively multiplying the amount surviving at the beginning of the age interval by one minus the interpolated or extrapolated retirement ratio for the age interval to obtain the amount surviving at the end of that age interval.

$$\text{Amount surviving at age } X + 1 = \left(\frac{\text{Amount surviving at age } X}{\text{Retirement ratio for age interval } X \text{ to } X + 1} \right) (1 - \text{Retirement ratio for age interval } X \text{ to } X + 1)$$

The linear regression equation is (16, pp. 382-383)

$$E(y|x_j) = \alpha + \beta x_j$$

y = dependent variable

x_j = independent variable

α, β = parameters

which is estimated by

$$E(y|x_j) = a + bx_j$$

$$y_j = a + bx_j + e_j$$

a = estimate of α

b = estimate of β

y_j = j^{th} value of the dependent variable

x_j = j^{th} value of the independent variable

e_j = deviation of the j^{th} observed value from the expected value of the j^{th} observation given that x_j is the independent variable

In the situation of fitting a polynomial to the retirement ratios

x_j = j^{th} age interval

y_j = observed retirement ratio for the j^{th} age interval

The unweighted least-squares method of fitting a linear regression line yields linear unbiased estimators of α and β , which have the minimum variance amongst the class of all linear unbiased estimators, if the following assumptions can reasonably be made (16, pp. 382-384):

1. The x_j values are controlled and/or measured without error.
2. The regression of y on x is linear, that is, $E(y|x_j) = \alpha + \beta x_j$.

3. The deviations $y_j - E(y|x_j)$ are mutually independent.
4. The deviations have the same variance (σ^2 , not usually known exactly) whatever be the value of x_j .

A fifth assumption is sometimes needed (16, p. 384):

In order to apply many standard statistical techniques, the further assumption that the conditional distribution of y , given x , is normal is needed.

This means the deviations mentioned in assumptions three and four, above, must be assumed to be normally distributed if ". . . many standard statistical techniques . . ." are to be used. Also, if this assumption (i.e., $e_j \sim N(0, \sigma^2)$) is valid, the least-squares method will yield unbiased estimators having the minimum variance amongst the class of all unbiased estimators (9, pp. 113-114).

The multiple linear regression model is of the form (16, p. 413)

$$E(y|x_1, x_2, \dots, x_k) = \alpha + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

where the x_i may be powers of the observed x 's, such as x_2 may be x_1^2 , x_3 may be x_1^3 , etc. The assumptions of the multiple linear regression model are similar to those of the simple linear regression model listed above.

The fourth assumption, above, is often called the assumption of homoscedasticity. If the assumption of homoscedasticity is invalid and if the variances are not known quantities, the least-squares estimators of the polynomial coefficients can be shown to be unbiased only under certain conditions; very little can be said about the variance properties of these estimators (9, p. 410).

The following example is presented to illustrate the plausibility of the non-constant variance of the retirement ratios from age interval to age interval (i.e., that the variance of the deviation $y_j - E(y|x_j)$ is not a constant, where y_j is the observed value of the retirement ratio at age

interval x_j). Fabricated data for each of three vintage groups are shown in Tables 1, 2, and 3. A composite of the retirement experiences of all three vintage groups is shown in Table 4. The three vintage groups are assumed to be three samples each of size 100 from the same parent population of property; the only difference between the units, as they are put into service, is the year of installation.

Table 1. Fabricated data for vintage group I

Age, years	Age interval, years	Age interval index number	No. surviving at beginning of age interval, units	No. retired during age interval, units	Retire- ment ratio
0	0-1	1	100	10	0.100
1	1-2	2	90	15	0.167
2	2-3	3	75	25	0.333
3	3-4	4	50	25	0.500
4	4-5	5	25	15	0.600
5	5-6	6	10	10	1.000
6	6-7	7	0	--	

Table 2. Fabricated data for vintage group II

Age, years	Age interval, years	Age interval index number	No. surviving at beginning of age interval, units	No. retired during age interval, units	Retire- ment ratio
0	0-1	1	100	8	0.080
1	1-2	2	92	15	0.163
2	2-3	3	77	23	0.299
3	3-4	4	54	29	0.537
4	4-5	5	25	13	0.520
5	5-6	6	12	10	0.833

Table 2 (Continued)

Age, years	Age interval, years	Age interval index number	No. surviving at beginning of age interval, units	No. retired during age interval, units	Retire- ment ratio
6	6-7	7	2	2	1.000
7	7-8	8	0	--	

Table 3. Fabricated data for vintage group III

Age, years	Age interval, years	Age interval index number	No. surviving at beginning of age interval, units	No. retired during age interval, units	Retire- ment ratio
0	0-1	1	100	13	0.130
1	1-2	2	87	18	0.207
2	2-3	3	69	29	0.420
3	3-4	4	40	28	0.700
4	4-5	5	12	12	1.000
5	5-6	6	0	--	

Table 4. Composite retirement experience of all three vintage groups

Age, years	Age interval, years	Age interval index number	No. surviving at beginning of age interval, units	No. retired during age interval, units	Retire- ment ratio
0	0-1	1	300	31	0.103
1	1-2	2	269	48	0.178
2	2-3	3	221	77	0.348
3	3-4	4	144	82	0.569
4	4-5	5	62	40	0.645
5	5-6	6	22	20	0.909

Table 4 (Continued)

Age, years	Age interval, years	Age interval index number	No. surviving at beginning of age interval, units	No. retired during age interval, units	Retire- ment ratio
6	6-7	7	2	2	1.000
7	7-8	8	0	--	

The graphs of the retirement ratios versus the age interval index numbers are shown in Figures 1, 2, and 3. Figure 4 shows the retirement ratios for all three vintage groups plotted on the same graph.

Several characteristics of retirement ratios should, perhaps, be noted:

1. Retirement ratios must be equal to or greater than zero and equal to or less than one.
2. The retirement ratio values of a vintage group generally increase, although not necessarily monotonically, as the age interval index number (age interval) increases.
3. A retirement ratio of one occurs only when all of the property of the vintage group surviving at the beginning of the age interval is retired during the age interval.
4. The observed retirement ratio at a given age interval is dependent, to some extent, upon the retirement ratios for all preceding age intervals because the number of units of a vintage group surviving at the beginning of an age interval (the denominator of the retirement ratio) is the number of units originally installed

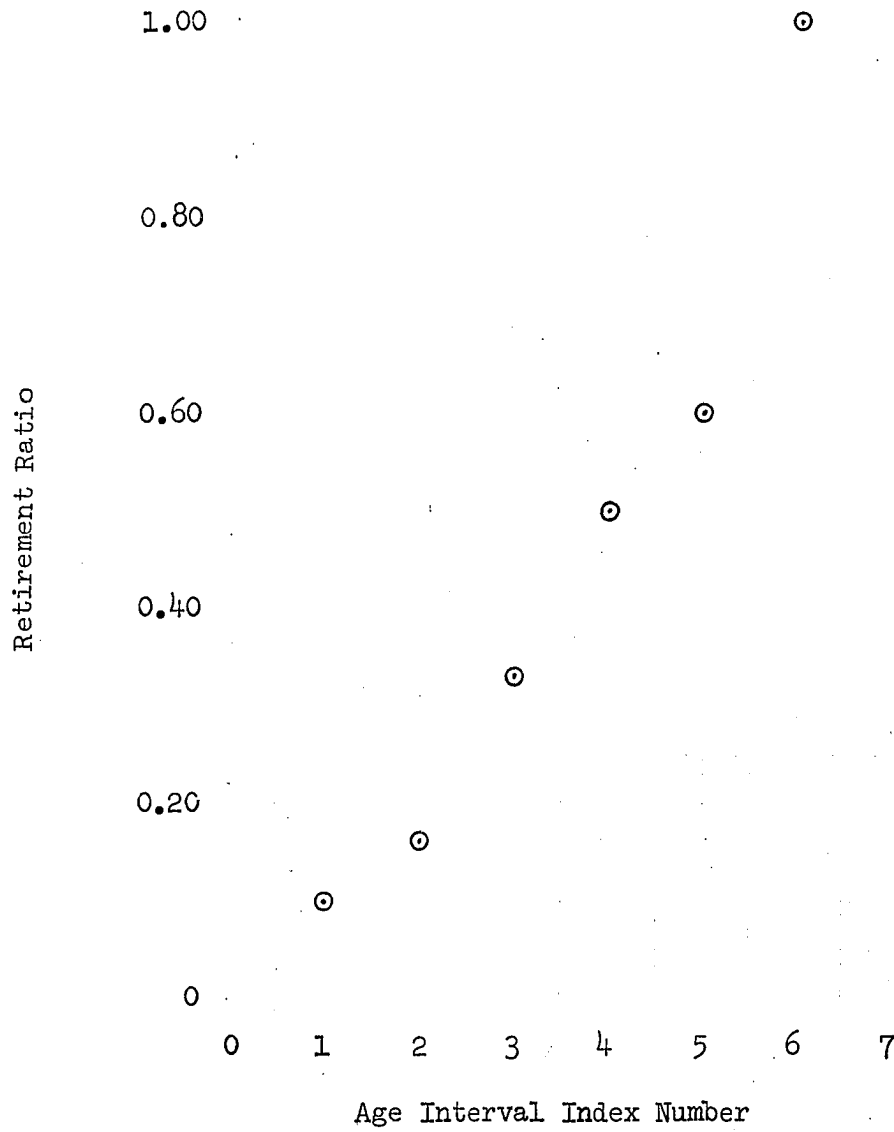


Figure 1. Retirement ratios for vintage group I

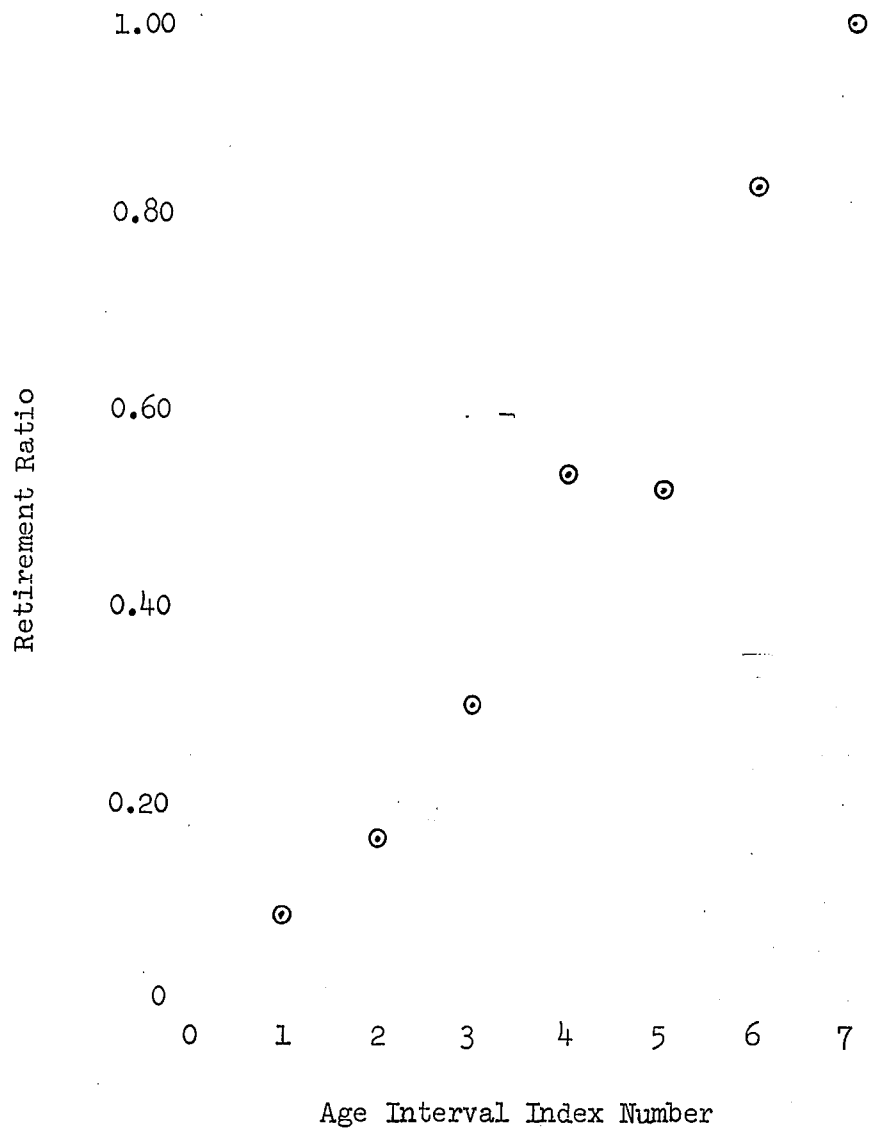


Figure 2. Retirement ratios for vintage group II

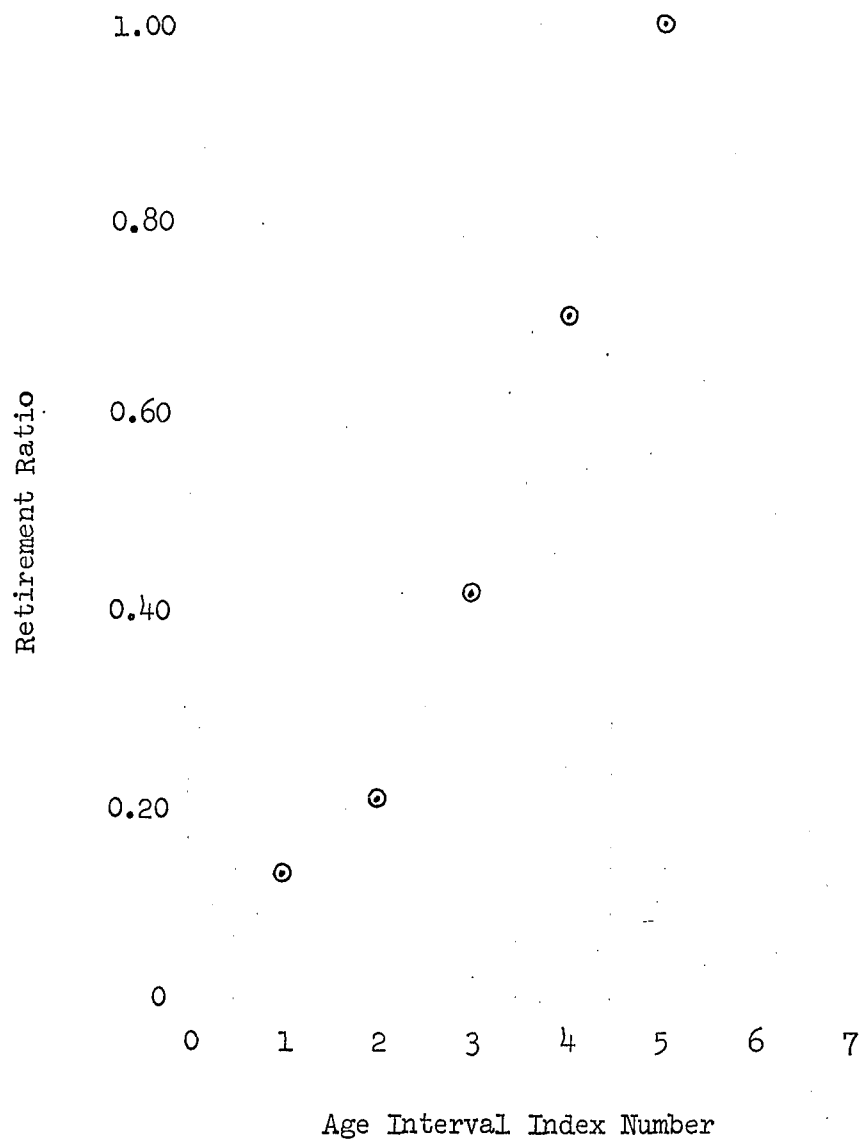


Figure 3. Retirement ratios for vintage group III



Figure 4. Retirement ratios of all three vintage groups

○ - vintage group I
□ - vintage group II
△ - vintage group III

less the numbers of units retired during preceding age intervals (the numerators of the retirement ratios of the preceding age intervals).

5. The several retirement ratios at each age interval (one from each vintage group) are not necessarily identical (see Figure 4), thus suggesting the possibility of a vertical dispersion or distribution of retirement ratios at each age interval.

The percent change in the retirement ratio, for a given change in the amount of property retired during the age interval, is dependent upon the amount of property surviving at the beginning of the age interval (the denominator of the retirement ratio). If

d_k = denominator of the retirement ratio for the k^{th} age interval
then

$$d_1 \geq d_2 \geq d_3 \geq \dots \geq d_K \quad k = 1, 2, \dots, K$$

Therefore, a given change in the amount of property retired during an age interval will generally cause a larger relative change in a retirement ratio if the retirement ratio is for a later age interval than if the retirement ratio is for an earlier age interval.

Figure 5 is the usual retirement ratio plot of the combined experience of all three vintage groups. The retirement ratio for the k^{th} age interval is calculated as the sum of the retirements from all vintage groups during the k^{th} age interval divided by the sum of the amounts surviving from all vintage groups at the beginning of the k^{th} age interval. Thus, Figure 5 illustrates a "horizontal" dispersion of retirement ratios across age intervals. The possibility of a "vertical" dispersion of retirement

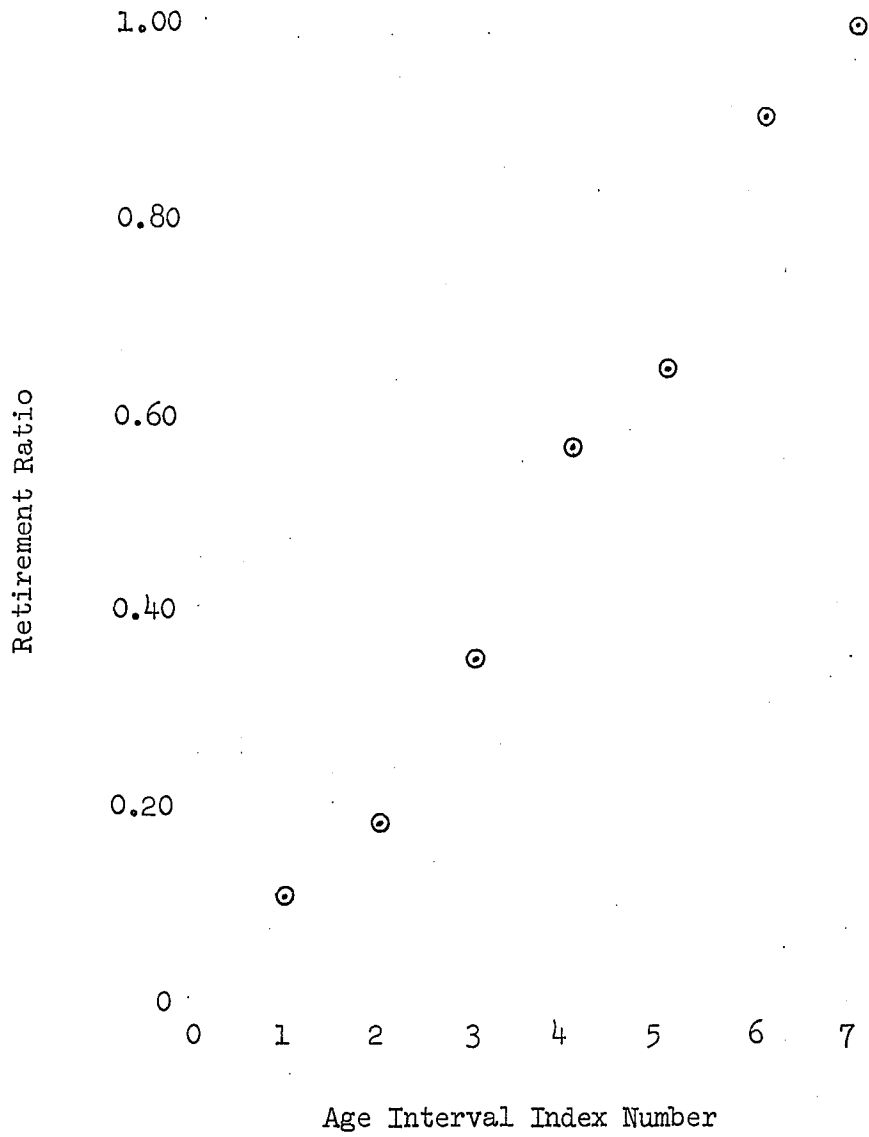


Figure 5. Composite retirement ratios

ratios within each age interval, in addition to the horizontal dispersion of retirement ratios across age intervals, is indicated in Figure 4.

Figure 4 also points up the possibility that the variance of the vertical distribution of retirements within an age interval is not necessarily the same as the variance of the vertical distribution of the retirement ratios within some other age interval.

The fourth assumption of the unweighted, least-squares method is (16, p. 383; see p. 9 of this dissertation):

These deviations (i.e., $y_j - E(y|x_j)$) have the same variance . . . whatever be the value of x_j .

As applied to fitting a polynomial to the retirement ratios, the above assumption means that the variance of the vertical distribution of retirement ratios at each age interval is assumed to be a constant (i.e., to be the same from age interval to age interval).

The subject of this investigation is the possible non-constant variance of the vertical distribution of retirement ratios and the effect of such on the method of fitting a polynomial to the retirement ratios.

OBJECTIVES OF INVESTIGATION

The smoothing, extending, and interpolating or extrapolating of retirement ratios is a method frequently used in life analysis to obtain a smoothed, complete life table or survivor curve. A number of assumptions (16, pp. 382-384; see pp. 8-9 of this dissertation) must be made if the unweighted, least-squares method of fitting a polynomial to the retirement ratios is to yield linear unbiased estimators of the polynomial coefficients having the minimum variance amongst the class of all linear unbiased estimators. Two of these assumptions, the third and the fourth, may not be valid. The third assumption (16, p. 383) appears to be invalid in view of the fact that the denominators of the retirement ratios, except the denominator of the retirement ratio for the first age interval, are dependent upon the preceding numerators. However, this third assumption is not a subject of investigation in this dissertation.

The subject of this investigation is the validity of the fourth assumption (of homoscedasticity) and a better means of fitting a polynomial to retirement ratios than the presently used, least-squares procedures if this assumption is invalid. A solution to the problems engendered by the failure of assumption four is dependent upon ascertaining the vertical distribution of the retirement ratios at each age interval and estimating certain parameters of these vertical distributions.

The specific objectives of this dissertation are:

1. To investigate the vertical distribution of retirement ratios at each age interval.

2. To investigate methods of obtaining estimators of certain parameters of the vertical distribution of retirement ratios at each age interval.
3. To develop, if possible, a more exact method of fitting a polynomial to the retirement ratios based on the findings in (1) and (2) above.

PRESENT ACTUARIAL METHODS OF LIFE ANALYSIS

A life analysis provides information for life estimation. The usefulness of the information is dependent mainly upon the appropriateness and reliability of the property data analyzed, the models used in analyzing the data, and the interpretation of the results. At best, the results of a life analysis provide more or less accurate estimates of the past mortality characteristics of the property in question. These results should be used in life estimation only to the extent that the past mortality behavior of the property is expected to be similar to the future mortality behavior of the property.

Two basic assumptions of life analysis, regardless of the methods used, are:

1. The mortality behavior of a property follows some "law of mortality" expressible in terms of time or some other variable.
2. The past mortality behavior of a property is indicative, to a greater or lesser extent, of the expected future mortality behavior of the property.

Although the "law of mortality" is generally expressed in terms of time, it could be expressed in terms of units of production or some other suitable variable(s). The first assumption implies that no extraneous variables make the relationship between retirements and time (or other variable) of little consequence. The extent to which the second assumption is incorrect is usually considered in the life estimation process.

Several methods of life analysis are called actuarial methods because of their similarity to methods developed by life insurance actuaries to

study human mortality (8, p. 44). The process of life analysis utilizing an actuarial method(s) can be roughly divided into three broad steps:

1. The selection of the property data to be analyzed and the aggregation of this data in a useful form.
2. The selection of the method(s) of obtaining the original life table or original survivor curve (or a derived curve, such as a retirement ratio curve), the selection of the particular data set(s) to which the method(s) is to be applied, and the application of the method(s) to the data set(s).
3. The selection and application of a method(s) of smoothing and a method(s) of extending (if necessary) the original life table (or original survivor curve or some other curve), and the interpolation and/or extrapolation of values to obtain a complete, smoothed life table.

In a life analysis study of a property, different methods of obtaining an original life table, etc. may be applied to different data sets and the resulting original life tables smoothed and extended by one or more methods to provide information about trends in mortality behavior.

Related Concepts

The actuarial methods of life analysis are based on statistical concepts. A basic concept of statistics is a probability distribution.

A table of the possible values which a chance event may assume with a corresponding probability for each value is called a probability distribution for the parent population (1, p. 19).

A mathematical function representing a probability distribution is called a probability distribution function (distribution function). The

probability of obtaining some value for an event which is equal to or less than a specified value is called the cumulative distribution; a mathematical function representing the cumulative distribution is called a cumulative distribution function. Distribution functions, and the corresponding cumulative distribution functions, may be either discrete or continuous functions; the appropriate form is dependent on whether the values the chance variate can take on are discrete or continuous.

The requirements of a function to be a discrete probability function are (24, p. 33)

$$1. f(x_i) \geq 0$$

$$2. \sum_{i=1}^N f(x_i) = 1 \quad i = 1, 2, \dots, N$$

where

x_i = the possible values which the chance variate, x , may assume

$f(x_i)$ = the probability that x takes on the values x_i ;

$$i = 1, 2, \dots, N$$

Then the cumulative distribution function, $F(a)$, is

$$\begin{aligned} F(a) &= \Pr(x \leq a) \\ &= \sum_{i=1}^a f(x_i) \end{aligned}$$

where "a" is some specified value of x . Also

$$\begin{aligned} F(a, b) &= \Pr(a \leq x \leq b) \\ &= \sum_{i=a}^b f(x_i) \end{aligned}$$

For the continuous distribution case, the assumptions are (24, p. 33)

$$1. f(x) \geq 0$$

$$2. \int_{-\infty}^{\infty} f(x)dx = 1$$

Then

$$\begin{aligned} F(a) &= \Pr(x \leq a) \\ &= \int_{-\infty}^a f(x) dx \\ F(a, b) &= \Pr(a \leq x \leq b) \\ &= \int_a^b f(x) dx \end{aligned}$$

The mathematical expectation of x , denoted as $E(x)$, is the mean or average value of x . $E(x)$ is the sum of the products of the distance of each x_i from the origin times the probability that the x_i will occur (i.e., the first moment of x about the origin). For the discrete case

$$E(x) = \sum_{i=1}^N x_i f(x_i)$$

and for the continuous case

$$E(x) = \int_{-\infty}^{\infty} x f(x) dx$$

The variance of x , denoted as σ^2 , is the second moment of x about the mean, $E(x)$. Let

$$\mu = E(x)$$

Then

$$\sigma^2 = E(x - \mu)^2$$

which for the discrete case is

$$\sigma^2 = \sum_{i=1}^N (x_i - \mu)^2 f(x_i)$$

and for the continuous case is

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx$$

If the individual service lives (ages at retirement) of the units comprising a property group are represented by x'_k , then $f(x')$ is the distribution function of the service lives. Service life could be measured on a discrete scale or on a continuous scale. The mathematics

of both the continuous case and the discrete case are presented because both the concept of a continuous scale and the concept of a discrete scale have proven to be useful. The presence of a summation symbol, Σ , in an equation indicates that service life is being considered as a discrete variable and the presence of an integral symbol, \int , indicates that service life is being considered as a continuous variable.

Even complete property accounting records generally show only the number of units in a vintage group and the year of installation of the vintage group and not the exact time that each unit of the vintage group was installed (and similarly for the time at which units are retired). For the discrete case, a common set of assumptions (often referred to as the half-year convention) is (20, pp. 147-148):

The assumption is made that the installations of a given calendar year were made somewhat uniformly throughout the year; therefore, the assumption that all the units were zero years old on July 1 of the year of installation is appropriate. The average age of retirements would then always be the integral years 1, 2, 3, etc. But retirements having an average age of, say 3 years, must be composed of units having specific ages varying from $2\frac{1}{2}$ to $3\frac{1}{2}$ years. Ages for specific reference in the calculation of the survivor curve or for a January 1 inventory date must be expressed on the $\frac{1}{2}$ -year basis.

Another customary assumption is that property retired during the same calendar year as it was installed is retired during the age interval $0-0\frac{1}{2}$, or at an average age of $0\frac{1}{4}$ year.

Therefore, for the discrete case only, let

x'_k = age index number

x'_1 = age $\frac{1}{4}$ years

x'_2 = age 1 year

.....

x'_K = age K-1 years

= maximum age index number

x_k = age interval index number

x_1 = age interval 0 to 1/2 years

x_2 = age interval 1/2 to 1 1/2 years

.....

x_K = age interval $(K - 1 \frac{1}{2})$ to $(K - \frac{1}{2})$ years

$f(x')$ = distribution function of the service lives

$f(x'_k)$ = probability of any unit being retired at age x'_k

$f(x)$ = distribution function, by age intervals, of the service lives

$f(x_k)$ = probability of any unit being retired during age interval k

Also, let

ω = maximum life when measured on a continuous scale

The average service life of a property group, ASL, is defined as the average age of all units at retirement (Figure 6).

$$\begin{aligned} \text{ASL} &= E(x) \\ &= \sum_{k=1}^K x_k f(x_k) \quad k = 1, 2, \dots, K \end{aligned} \quad (1)$$

where

$$\sum_{k=1}^K f(x_k) = 1$$

Also

$$\text{ASL} = \int_0^{\omega} x f(x) dx \quad 0 \leq x \leq \omega$$

The cumulative distribution function (Figure 7) is

$$F(a) = \sum_{k=1}^a f(x_k)$$

or

$$F(a) = \int_0^a f(x) dx$$

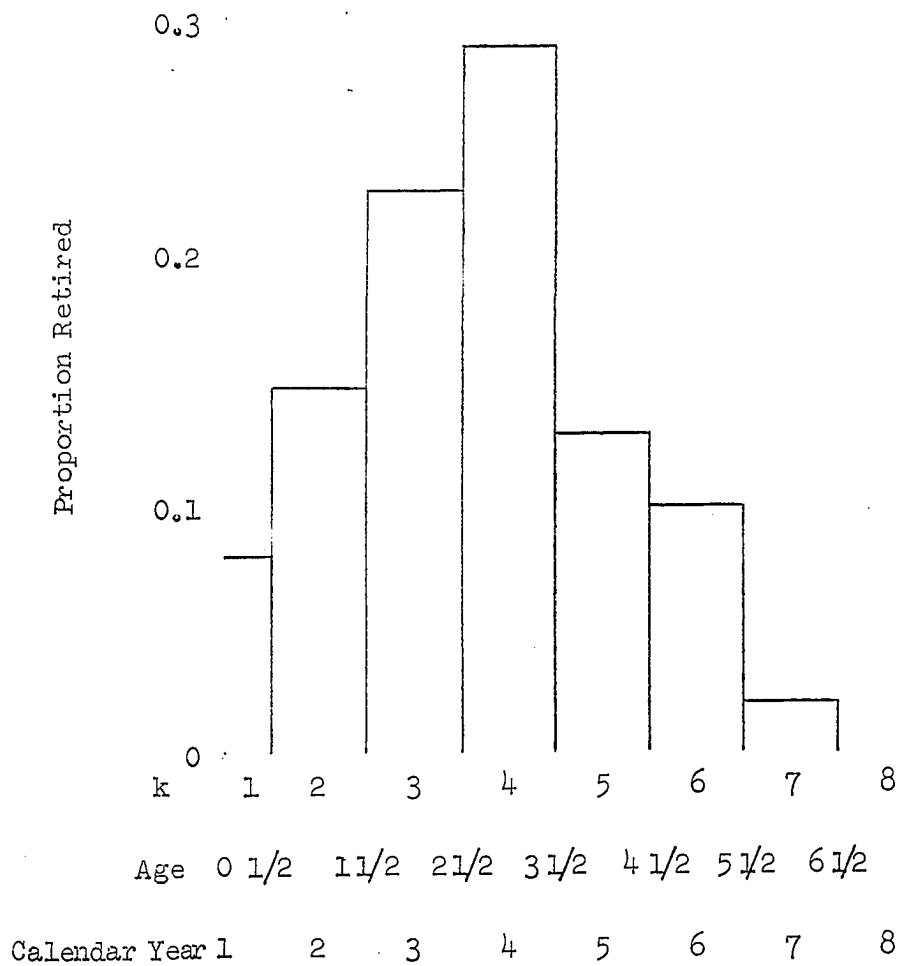


Figure 6. Frequency distribution of retirements

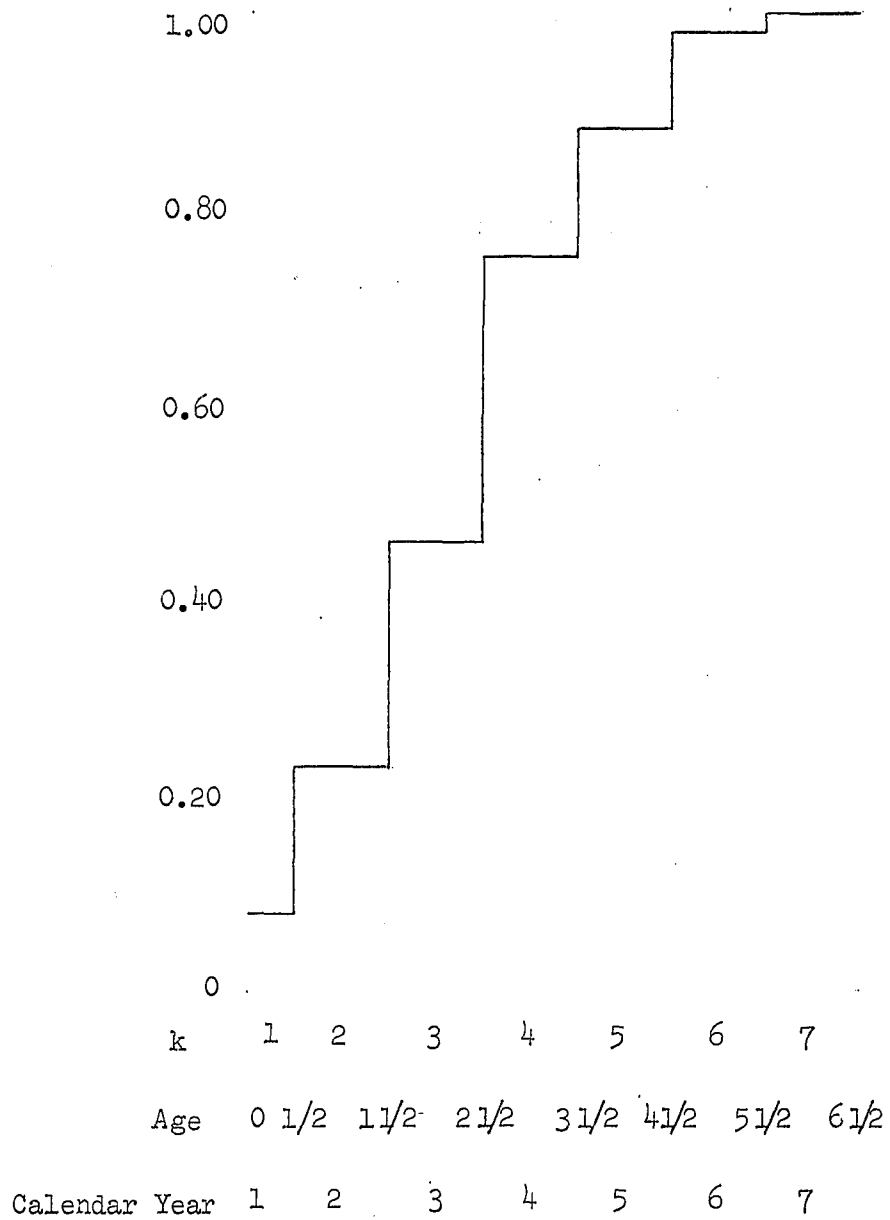


Figure 7. Cumulative distribution of retirements

where

a = some specified age

The survivor function or survivor curve, $y(x)$, represents the proportion of units surviving at any age (Figure 8).

$$y(a - \frac{1}{2}) = 1 - \sum_{k=1}^a f(x'_k)$$

or

$$y(a) = 1 - \int_0^a f(x)dx$$

Also, since

$$\sum_{k=1}^K f(x'_k) = 1$$

$$\int_0^{\omega} f(x)dx = 1$$

Then

$$y(a - \frac{1}{2}) = \sum_{k=a+1}^K f(x'_k)$$

$$y(a) = \int_a^{\omega} f(x)dx$$

For example

$$\begin{aligned} y(3 - \frac{1}{2}) &= y(2 \frac{1}{2}) \\ &= 1 - \sum_{k=1}^3 f(x'_k) \\ &= 1 - f(x'_1) - f(x'_2) - f(x'_3) \\ &= 1 - \text{retirements of average age } 1\frac{1}{4} - \text{retirements of} \\ &\quad \text{average age } 1 - \text{retirements of average age } 2 \\ &= \text{survivors at age } 2 \frac{1}{2} \\ &= \sum_{k=a+1}^K f(x'_k) \\ &= f(x'_4) + f(x'_5) + \dots + f(x'_K) \end{aligned}$$

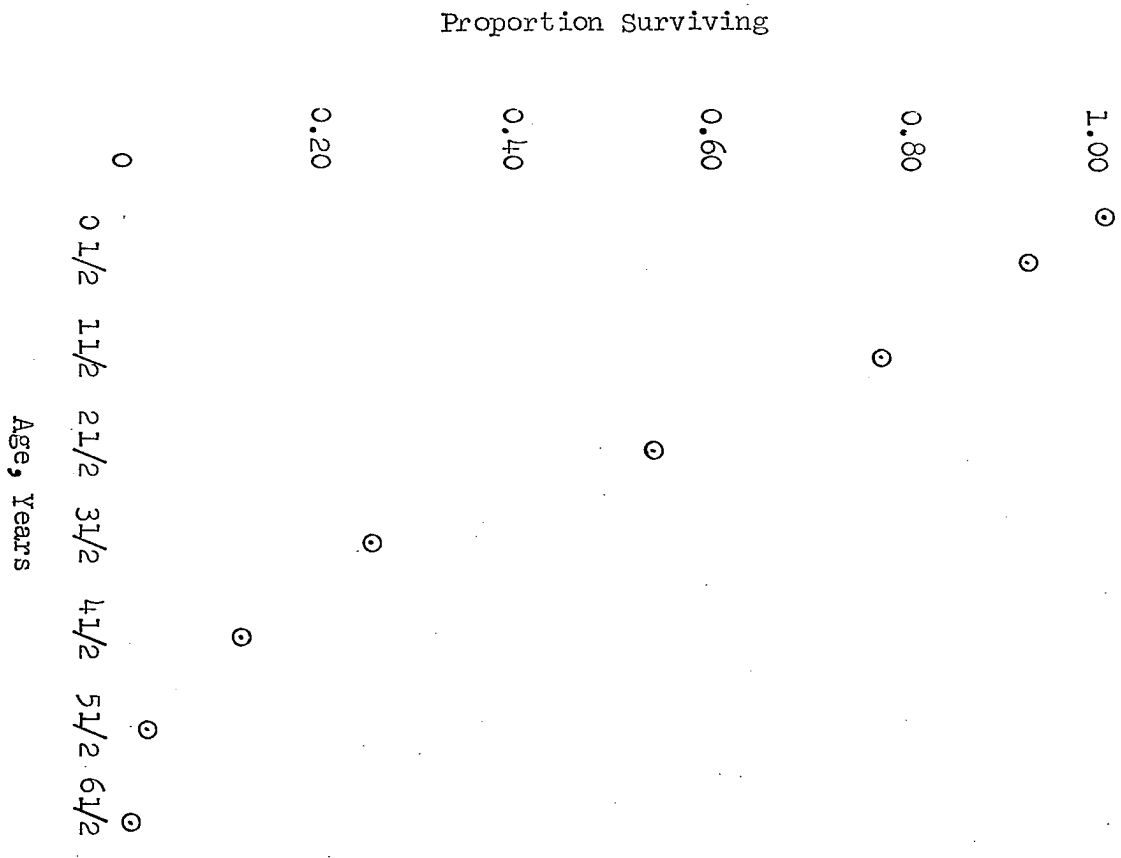


Figure 8. Survivor curve

= retirements of average age 3 + retirements of average age 4 + . . . + retirements of average age K - 1

The area under the survivor curve can be approximated by numerical integration and is equal to the average service life (as will be shown). In this particular case, finding the area under the curve using horizontal area strips is convenient. The difference in height of two successive points on the survivor curve (ages 0 and 0.5, ages 0.5 and 1.5, etc.) is the width of the horizontal area strip and is just $f(x_k)$. The average height of the horizontal area strip is the distance from the y axis to mid-way between the two points on the survivor curve (or one-half of the sum of the ages at the beginning and end of age interval k).

$$\begin{aligned} \text{ASL} &= (1/4) f(x_1) + (1) f(x_2) + \dots + (K - 1) f(x_K) \\ &= \sum_{k=1}^K x_k f(x_k) \end{aligned} \quad (2)$$

Equation 2 is the same as equation 1, thus showing that the area under the survivor curve is equal to the average service life. The average service life can also be calculated as the first moment of the frequency distribution about the origin.

$$\begin{aligned} \text{ASL} &= \sum_{k=1}^K x'_k f(x'_k) \\ &= (1/4) f(x'_1) + (1) f(x'_2) + \dots + (K - 1) f(x'_K) \end{aligned}$$

The retirement ratio, r_k , for an age interval k is defined as

$$\begin{aligned} r_k &= \frac{\text{number of units retired during age interval k}}{\text{number of units surviving at the beginning of age interval k}} \\ &= \frac{\text{proportion retired during age interval k}}{\text{proportion surviving at beginning of age interval k}} \end{aligned}$$

In the discrete case, $f(x_a)$ represents the proportion of the original

placement of units retired during the age interval a .

$$r_a = \frac{f(x_a)}{a-1} \cdot \frac{1}{1 - \sum_{k=1} f(x_k)}$$

$$= \frac{f(x_a)}{\sum_{k=a}^K f(x_k)}$$

The retirement ratio for the continuous case is

$$r_a = \frac{f(a)dx}{1 - \int_0^a f(x)dx}$$

$$= \frac{f(a)dx}{y(a)}$$

$$= \frac{-d[y(a)]}{y(a)}$$

The survival ratio, s_a , for an age interval a is defined as

$$s_a = \frac{\text{number of units surviving at end of age interval } a}{\text{number of units surviving at beginning of age interval } a}$$

$$= \frac{\text{proportion of units surviving at end of age interval } a}{\text{proportion of units surviving at beginning of age interval } a}$$

$$= \frac{\sum_{k=a+1}^K f(x_k)}{\sum_{k=a}^K f(x_k)}$$

For the continuous case

$$s_a = \frac{y(a) + d[y(a)]}{y(a)}$$

since the number or proportion of units retired during the small interval of time after $y(a)$ is $-d[y(a)]$.

The survival ratio is also equal to one minus the retirement ratio.

$$\begin{aligned}
 s_a &= 1 - \frac{f(x_a)}{\sum_{k=a}^K f(x_k)} \\
 &= \frac{\sum_{k=a}^K f(x_k) - f(x_a)}{\sum_{k=a}^K f(x_k)} \\
 &= \frac{\sum_{k=a+1}^K f(x_k)}{\sum_{k=a}^K f(x_k)}
 \end{aligned}$$

or

$$\begin{aligned}
 s_a &= 1 - \left\{ - \frac{d[y(a)]}{y(a)} \right\} \\
 &= \frac{y(a) + d[y(a)]}{y(a)}
 \end{aligned}$$

Expectancy is defined as (28, p. 12):

. . . that period of time extending from the observation age (usually the present) to the average of the forecasted dates when the units probably will be retired.

Expectancy is the future average years of service expected from the units surviving at the observation time. The expectancy at any age, $E_a - 1/2$, is

$$\begin{aligned}
 E_a - 1/2 &= \frac{\text{area under the survivor curve to the right of age } a - 1/2}{\text{proportion surviving at age } a - 1/2} \\
 &= \frac{\sum_{k=a+1}^K (k - a - 1/2) f(x_k)}{\sum_{k=a+1}^K f(x_k)}
 \end{aligned}$$

or

$$E_a = \frac{\int_a^{\omega} [1 - \int_0^a f(x)dx]dx}{1 - \int_0^a f(x)dx}$$

$$= \frac{\int_a^{\omega} y(x)dx}{y(a)}$$

The probable life of the units surviving at any age, $P_a - 1/2$, is defined as

$$P_a - 1/2 = a - 1/2 + E_a - 1/2$$

If the frequency with which x'_k occurs (rather than the proportion retired at an average age of x'_k) is known, the distribution function, cumulative distribution function, average service life, etc., may be stated as follows for the discrete case. Let

$f(x''_k)$ = frequency with which x'_k occurs

x'_k = average age at retirement (includes the ages $x'_k - 1/2$ to $x'_k + 1/2$)

Then

$$f(x'_k) = \frac{f(x''_k)}{K \sum_{k=1} f(x''_k)}$$

= proportion retired at an average age of x'_k ; the distribution function

$$F(a) = \frac{\sum_{k=1}^a f(x''_k)}{K \sum_{k=1} f(x''_k)}$$

$$ASL = \frac{\sum_{k=1}^K x_k' f(x_k'')}{\sum_{k=1}^K f(x_k'')}$$

The number of units surviving at any age is

$$y'(a - 1/2) = \sum_{k=1}^K f(x_k'') - \sum_{k=1}^a f(x_k'')$$

The proportion surviving at any age is

$$\begin{aligned} y(a - 1/2) &= \frac{\sum_{k=1}^K f(x_k'') - \sum_{k=1}^a f(x_k'')}{\sum_{k=1}^K f(x_k'')} \\ &= 1 - \frac{\sum_{k=1}^a f(x_k'')}{\sum_{k=1}^K f(x_k'')} \end{aligned}$$

Also

$$r_a = \frac{\frac{f(x_a'')}{\sum_{k=1}^K f(x_k'')}}{\frac{\sum_{k=a}^K f(x_k'')}{\sum_{k=1}^K f(x_k'')}} = \frac{f(x_a'')}{\sum_{k=a}^K f(x_k'')}$$

$$= \frac{f(x_a'')}{\sum_{k=a}^K f(x_k'')}$$

$$s_a = \frac{\sum_{k=a+1}^K f(x_k'') / \sum_{k=1}^K f(x_k'')}{\sum_{k=a}^K f(x_k'') / \sum_{k=1}^K f(x_k'')} = \frac{\sum_{k=a+1}^K f(x_k'')}{\sum_{k=a}^K f(x_k'')}$$

$$= \frac{\sum_{k=a+1}^K f(x_k'')}{\sum_{k=a}^K f(x_k'')}$$

$$\begin{aligned}
 E_a - 1/2 &= \frac{\sum_{k=a+1}^K \{(k - a - 1/2)[f(x_k'') / \sum_{k=1}^K f(x_k'')]\}}{\sum_{k=a+1}^K [f(x_k'') / \sum_{k=1}^K f(x_k'')]} \\
 &= \frac{\sum_{k=a+1}^K (k - a - 1/2) f(x_k'')}{\sum_{k=a+1}^K f(x_k'')}
 \end{aligned}$$

$$P_a - 1/2 = a - 1/2 + E_a - 1/2$$

The amount surviving at any age is often expressed as a percent. The percent surviving can be calculated by multiplying the proportion surviving by 100%.

Mathematical expressions for the continuous case can be derived in a similar manner.

Selection and Aggregation of Property Data

The data to analyze for the purpose of predicting the mortality behavior of a property are, generally, the historical data of that property. Certain assumptions are generally made about the property data:

1. Historical data on the same or a similar property group are available.
2. The property group is composed of homogeneous units or of different units in substantially the same relative amounts as are expected in the future.
3. Sufficient data in a usable form are available to make an actuarial life analysis.

Historical data on the property which is the subject of a life analysis may not be available or may be unusable. In this case, the analyst may analyze the historical data of a property which he thinks will exhibit mortality behavior similar to that of the property which is the subject of the life estimation process. The results of a study of a similar property should be given only such weight in the life estimation process as is appropriate. Another, infrequently used alternative is to take a complete inventory of the property. A third alternative is to proceed directly to the life estimation process without making a life analysis study; the analyst's knowledge and his experience in life analysis provides the type of information which is usually obtained through life analysis.

The subject of the homogeneity of a property, for life analysis purposes, involves at least two areas: (1) the physical characteristics of the property and (2) the measure of the amount of property. A property group account may include several different sizes and types of property. Even if only one size and type of property is recorded in a given account, heterogeneity may arise from including property manufactured in different years and which may be different because of modifications in materials and/or design.

Two common measures of the amount of property are physical units and dollars, the latter being the most frequently used measure. The age at retirement of one physical unit may be (and often is) independent of the age at retirement of any other physical unit. On the other hand, the physical units comprising a vintage group are often heterogeneous because of their different physical characteristics.

Dollars are homogeneous in the sense that one dollar is numerically equal to another dollar. Thus, dollars provide a common scale for measuring amount of property. However, the number of dollars invested in one item of a property group generally is not the same as the number of dollars invested in another item of the property group. The age at retirement of one dollar is rarely independent of the age at retirement of some other dollar(s). Hence, dollars are not independent random variables; a fact which might inhibit development of a statistical procedure for fitting polynomials to retirement ratios.

Howard (15) compared the average lives and the accrued depreciation, for group property, computed on (1) a unit basis and (2) a dollar basis. Data on the mortality experience of freight cars were used. Complete physical-unit data were available for the years 1918 through 1945. The only dollar data available were ". . . the total dollars remaining in service January 1 of each year, the total dollars placed in service each year, and the total dollars retired from service each year" (15, p. 19). Therefore an average unit cost for each unit installed in a vintage year was calculated by dividing the total dollars placed in service that year by the number of freight cars placed in service that year. The results of his study indicated that the average service lives of the freight cars were, at most, one-half year greater when calculated on the dollar basis than when calculated on the unit basis. The difference in average service lives on a unit basis and on a dollar basis was attributed to the greater weight given to more recent placements because of a rising price level.

The results are indicative of the differences to be expected because of price level changes. They are not indicative of the differences to be expected because of between-unit price differences because of his averaging of unit costs within a year.

The records of the property group selected for study must be carefully reviewed. The data may need to be adjusted for a number of reasons, such as: the type(s) of property included in the property group has been changed from time to time or accounting practices have been changed from time to time or data on properties which have been sold or acquired as used, but useful, property have been included in the record (5, pp. 7-10).

Despite the importance of this factor, and the fact that the time required to correct and adjust the books of account is ordinarily many times the manhours required to make the statistical analyses themselves, the literature on this subject is not very helpful. It is replete with warnings that early book records are often incomplete, that accounting distinctions between capital and maintenance charges have undergone changes, that the type of equipment represented by a given plant account may change from one generation to the next, and so forth, but it contains very little by way of specific suggestions for approved treatment of the raw data to make them suitable for analysis (5, p. 10).

Original Life Table

An original life table is a tabulation from the raw data of the amount of property surviving from an original placement at each age. A plot of the amount surviving versus age (generally on rectangular coordinate graph paper) is called an original survivor curve.

The amounts surviving may be expressed as physical units, dollars, proportions, or percents. If the original life table or original survivor curve is expressed in terms of percents or proportions, a more direct comparison can be made between different life tables or survivor curves.

Percent surviving is commonly used and will be used hereafter unless otherwise noted.

The original life table may be obtained from the raw data in at least five ways: individual-unit method, original-group method, composite original-group method, multiple original-group method, and annual-rate method (28, pp. 17-18). The choice of method(s) to use is dependent on the data available, the purpose of the life analysis, and the type of information to be obtained in applying the method. As mentioned previously, one or more methods may be applied to several different data sets to obtain information about the mortality behavior of the property.

The original-group, composite original-group, and annual rate methods require relatively complete historical data covering the years of experience of the vintage group(s) included in the analysis. The individual-unit and multiple original-group methods require less complete data but yield less useful results than the other three methods. These two methods are used, generally, only if the data available are insufficient to permit use of any of the other three methods.

The results obtained by the application of these five methods to the historical data of property will usually be different. These differences in results are due to various factors, such as:

1. Use of different data sets,
2. Random variation in sample data, and
3. Changes in the mortality behavior exhibited by the property group resulting from changes in those factors influencing the retirement of property.

Individual-unit method

The individual-unit method can be used when the only data available are the amounts and ages at retirement of property retired during a calendar year or several adjacent calendar years. The retirements during the calendar year(s) are arranged in ascending order according to age at retirement. The sum of all such retirements is taken to be the total amount of property "surviving" at age zero. The percent surviving at each successive age or the percent surviving at the beginning of each successive age interval is the amount of the retired property that was retired at a later age. The original life table will always extend to zero percent surviving because only retired property is considered in calculating the table.

The average service life obtained by numerical integration is the average age at retirement of those units retired during the calendar year(s) of observation, not the probable average service life of the property. If the property has not reached stability (i.e., no growth, no decline, and renewals approximately equal to retirements), the average age at retirement may not be a very good approximation of the probable average service life of the property. Similarly, the retirement dispersion pattern obtained may not be a very good approximation of the probable retirement dispersion pattern of the property.

Original-group method

Data required for the original-group method are:

1. The amount of the property installed in a given year, a vintage group, and

2. The amounts of and ages at retirement of the property already retired.

If the original life table is incomplete, the table or the corresponding survivor curve will have to be extended to zero percent surviving before the mortality characteristics can be ascertained.

The calculated, probable average service life is the probable average service life of the particular vintage group. A study of successive vintage groups may indicate trends over time, if any, of the mortality characteristics of the property due to changes in the physical characteristics of the property.

Composite original-group method

The composite original-group method treats the combined mortality experience of two or more vintage groups as the mortality experience of a single group. Data requirements are similar to those of the original-group method. If an incomplete, original life table is obtained, the table or the survivor curve must be extended to zero percent surviving before the mortality characteristics can be ascertained.

This method is especially useful when only a relatively small amount of property is installed each year and/or the mortality experience of a single vintage group is erratic. The mortality characteristics obtained are composites of the mortality characteristics of the individual vintage groups included in the single combined group.

A rolling-band study, a series of analyses of different composite groups, may indicate trends in the mortality characteristics of the property over time. Each successive composite group is formed from the

preceding composite group by eliminating the oldest (or youngest) vintage group in the composite group and adding the vintage group just subsequent to (or just preceding) the composite group.

As the number of vintage groups included in the composite group increases, the mortality experience of the composite group tends to become less erratic. On the other hand, grouping a large number of vintage groups into a single group tends to mask trends in the mortality behavior of the property.

Multiple original-group method

The multiple original-group method requires data on the ages and amounts of the property surviving as of a given date. A table of the ages and amounts surviving, arranged in order of increasing age, constitutes the original life table. Percent surviving values can be calculated by using the amount surviving from the most recent vintage group as the denominator of the fraction

$$\text{percent surviving at age } x = \frac{\text{amount surviving at age } x}{\text{amount surviving at age zero}} (100\%)$$

If the percent surviving at any age exceeds 100%, when calculated in the above manner, the common practice is to reduce such values to 100%.

Successive entries in the original life table may be larger or smaller than previous or subsequent entries because:

1. Each vintage group provides one entry in the table,
2. The amount of property surviving from a vintage group is related to the amount installed during that vintage year, and
3. No consideration is given to the various amounts installed during each vintage year nor to the amounts already retired from each

vintage group.

Unless the property has reached stability, the original life table and original survivor curve tend to be erratic and incomplete.

Annual-rate method

By this method, the original life table is calculated from that mortality experience of a number of vintage groups (called the placement band) exhibited during a given period of years (called the observation band). The data required on each vintage group included in the placement band are:

1. The ages and amounts of property (in units or dollars) retired each year during the observation band of years and
2. The amount of property surviving at the beginning of each year that the vintage group is included in the observation band.

A retirement ratio for each age interval is calculated as follows

$$r_{x - 1/2 \text{ to } x + 1/2} = \frac{\sum_{i=1}^I \text{property from the } i^{\text{th}} \text{ vintage group retired during the age interval } x - 1/2 \text{ to } x + 1/2 \text{ during the observation band of years}}{\sum_{i=1}^I \text{property from the } i^{\text{th}} \text{ vintage group surviving at age } x - 1/2 \text{ during the observation band of years}}$$

i = index number of the vintage group

= 1, 2, . . . , I

The percent surviving at the end of each age interval can then be calculated by starting with 100% surviving at age zero and successively multiplying the percent surviving at the beginning of each age interval by one minus the retirement ratio for that age interval.

If sufficient data are available, the annual-rate method is generally one of the methods used to obtain original life tables in a life analysis because:

1. The mortality experience of the most recent vintage years can be utilized,
2. The mortality behavior of the property during the observation band of years reflects the effects of management policies, economic conditions, public requirements, etc., on the retirement of property during the observation band of years, and
3. Both property surviving and property retired are considered.

A rolling band type of analysis, in which the most recent (or earliest) year of the observation band is eliminated and the year preceding (or subsequent to) the earliest (or most recent) year of the observation band is added, is frequently made to study any trends in the mortality behavior of the property.

Marston et al. (20, p. 154) suggest an observation band of three to thirty years. An observation band of only a few years permits the more recent mortality experience of the property to exert a greater influence on the values in the original life table. On the other hand, an original life table based on a narrow observation band is more likely to be erratic than an original life table based on a relatively wide observation band.

Methods of Obtaining a Smoothed Life Table

The original life table or original survivor curve is frequently incomplete because not all of the units of even the oldest vintage groups included in the data set have been retired. A complete life table or

survivor curve must be obtained before the mortality characteristics of the property can be ascertained.

The process of obtaining a complete survivor curve is composed of two steps: (1) fitting a smooth curve to the existing data and (2) extending the smoothed curve to zero percent surviving. A smooth curve may be fitted to the available data by a variety of methods, such as the various matching and mathematical methods. Extending the smoothed curve to zero percent surviving is a matter of judgment. Where the method of smoothing provides an "extension" of the curve, this extension is often accepted unless it is obviously incorrect. A more appropriate approach is to use judgment to select the most likely extension of the curve, the extension obtained from the smoothing step being considered as only one of the possible alternatives.

Three general methods of fitting a smooth curve to the raw data are judgment, matching to type curves, and statistical methods. Even if a complete, original life table is obtained, a smooth curve is often fitted to the data, by one of the above methods, before the mortality characteristics are ascertained. Marston et al., with reference to estimating the probable average service life from the survivor curve, say (20, p. 164):

The stub curve must be extended to zero percent surviving and the irregular curve should be smoothed before the average service life is computed. The objective is to obtain the most probable average service life. Such probability is indicated by a smooth complete survivor curve because such a smooth curve is the type most likely to result from observations at regular yearly intervals of large numbers of exposures to retirements.

Judgment method

Smoothing the survivor curve by judgment is accomplished by plotting the percent surviving at each age on rectangular coordinate graph paper and drawing, by judgment, a smooth curve through the points. Extension of the survivor curve to zero percent surviving is frequently accomplished by judgment, also. Obviously, no two analysts given the same set of points on a survivor curve are likely to draw exactly the same smooth curve; however, the difference between two such smoothed curves may be negligible from a practical point of view.

Numerical integration of the survivor curve yields the probable average service life. Additional calculations are required to obtain the expectancy and probable life at each age. Utilization of a high-speed digital computer would greatly reduce the time and effort involved in numerical integration and in subsequent calculations.

Matching method

The matching method involves comparing the original survivor curve, or a related curve, to a family of standardized curves and selecting that member of the family of curves which best fits or represents the data points. The criterion for determining which member of the family best fits the data is generally judgment. Other criteria may be used, such as selecting that member of the family of curves which minimizes the sum of the squares of the differences between the members of the family and the original survivor curve.

The Iowa type curves are the most widely recognized family of standard curves (4, p. 19). Bulletin 125 Revised, of the Iowa State University

Engineering Research Institute (28), contains all twenty-two of the Iowa type curves. The original eighteen Iowa type curves, developed by Winfrey and Kurtz (28, 29), are divided into three sets on the basis of the position of the mode of the frequency curve with respect to the average service life: six left-modal, seven symmetrical, and five right-modal. Couch (3) developed three origin-modal curves and, also, the data for the straight line survivor curve, in 1957. All four of these curves were designated as origin-modal.

A common procedure in using the Iowa curves is to first plot the original survivor curve points on transparent, rectangular coordinate graph paper. The standard Iowa curves are drawn on graph paper to a similar scale and for various average service lives. The plot of the original survivor curve points is superimposed on the graphs of the standard curves and the best fitting standard curve chosen by judgment. Winfrey suggests drawing a smooth curve, by eye, through the points of the original survivor curve before comparing the plot to the standard curves (28, p. 85).

Hoover (14) investigated the possibility of using an analog computer to match standard curves to the original survivor curve points. The circuitry for the following types of standard curves or functions were developed:

1. Iowa type survivor curves,
2. Weibull survivor function,
3. Gompertz-Makeham survivor function,
4. Truncated normal distribution function, and
5. Polynomial retirement ratio function.

Hoover matched the Iowa type curves to the stub data developed by Cowles (4). The standard curves were successively generated and displayed on an oscilloscope to which was taped a plot of the original survivor curve points; the curve type and average service life were controlled by the computer operator. The best fitting standard curve was selected by judgment. Hoover concluded that the analog computer could be used to develop estimates of mortality dispersion pattern and average service life.

Kimball (17) developed a family of curves (called the h-type curves) based upon a truncated normal distribution of retirements. The truncation of the frequency distribution occurs to the left of the average service life (i.e., at age zero). Hence, the h-type curves are generally left-modal. The relatively high-modal curves are essentially symmetrical and have small variance. As the modal value decreases, the variance increases and the frequency distribution becomes more and more left-modal with the negative exponential as the limiting form.

Although the h-system of life tables is of course not applicable to all cases of property retirements, for purposes of the general consideration of the behavior of property retirements in the broader aspects of the problem it is very useful to have such a system of life tables available in simple mathematical form. Tests of this system against several hundred life tables based on actual experience of utility property studied in the Bureau of Valuation of the New York Commission indicate very close agreement (17, p. 359).

Other families of type curves have been developed, such as the Patterson curves (22, pp. 60-68).

Statistical curve fitting methods

Statistical methods of fitting a smooth curve to raw data exhibit a number of desirable characteristics:

1. Any two analysts utilizing the same mathematical model and fitting technique to fit a smooth curve to the same raw data points should obtain the same results,
2. High-speed digital computers can be utilized making it possible to analyze a large number of data sets in a relatively short period of time, and
3. Human judgment is eliminated from the process of fitting a smoothed curve to the raw data points (this may, at times, be undesirable).

Judgment must be used in selecting the mathematical model, the fitting technique, and the data sets to be analyzed.

An extension of the smooth curve beyond the raw data points can be obtained from the mathematical function. Whether such an extension is reasonable or not is a matter of judgment. Unfortunately, mathematical methods foster an opposite approach, that of accepting the mathematical extension unless the extension is clearly unreasonable.

Mathematical functions can be fitted to a number of different, but related, sets of data points or ratios, such as the observed life table, the survival ratios, the retirement ratios, the retirement frequency distribution, or the cumulative retirements.

A frequently used method of obtaining a smoothed life table is the retirement ratio method. A function, such as a polynomial or power

function, is selected by judgment and the function fitted to the retirement ratios by the method of least-squares (5, p. 15). Polynomials are, perhaps, the most frequently used functions and will be used for illustrative purposes.

The retirement ratio at an age interval, say k , is defined as

$$r_k = \frac{f(x_k)}{\sum_k f(x_k)} \quad k = 1, 2, \dots, K$$

$f(x_k)$ = number of units retired during age interval k

If the composite-original group method or the annual-rate method is used to obtain the original life table, the experience of several vintage groups must be combined to obtain the retirement ratio at each age interval. A weighted average retirement ratio (rather than the average of the several retirement ratios) is usually calculated. Let

$$S_{ik} = \text{number of units surviving at the beginning of the } k^{\text{th}} \text{ age interval from the } i^{\text{th}} \text{ vintage group contributing experience to be included in the } k^{\text{th}} \text{ age interval retirement ratio}$$

$$= \sum_k f(x_{ik})$$

$$R_{ik} = \text{number of units retired during the } k^{\text{th}} \text{ age interval from the } i^{\text{th}} \text{ vintage group contributing experience to be included in the } k^{\text{th}} \text{ age interval retirement ratio}$$

$$= f(x_{ik})$$

Then

$$r_{ik} = \frac{R_{ik}}{S_{ik}}$$

The weight given the retirement ratio of each vintage group is the number of units of that vintage group exposed to retirement at the beginning of the age interval.

$$\begin{aligned} r_{\cdot k} &= \text{weighted average retirement ratio} \\ &= \frac{S_{1k} r_{1k} + S_{2k} r_{2k} + \dots + S_{Ik} r_{Ik}}{S_{1k} + S_{2k} + \dots + S_{Ik}} \end{aligned}$$

But

$$r_{1k} = \frac{R_{1k}}{S_{1k}}$$

therefore

$$\begin{aligned} r_{\cdot k} &= \frac{R_{1k} + R_{2k} + \dots + R_{Ik}}{S_{1k} + S_{2k} + \dots + S_{Ik}} \\ &= \frac{\sum_{i=1}^I R_{ik}}{\sum_{i=1}^I S_{ik}} \\ &= \frac{\sum_{i=1}^I f(x_{ik})}{\sum_{i=1}^I \sum_{k=1}^K f(x_{ik})} \end{aligned}$$

The function to minimize, when fitting a polynomial to these retirement ratios by the unweighted, least-squares fitting technique, is

$$\text{Min}_{a,b,c,\text{etc.}} \sum_{k=1}^K [r_{\cdot k} - (a + bx_k + cx_k^2 + \dots)]^2$$

Quite often a weighted least-squares fit is made by weighting the weighted average retirement ratio at each age interval by the total number

of units surviving at the beginning of that age interval.

Because the several plotted points do not carry equal weight, as pointed out before, it may be felt worth while to weight each according to the dollars or number of physical units involved (5, p. 15).

The function to minimize for the weighted least-squares fit is:

$$\text{Min}_{a,b,c,\text{etc.}} \sum_{k=1}^K \left\{ \sum_{i=1}^I S_{ik} [r_{.k} - (a + bx_k + cx_k^2 + \dots)]^2 \right\}$$

Only rarely is a polynomial of the fourth degree, or higher, selected as best representing the retirement ratio curve (21, p. 248).

Cowles (4) compared the results of smoothing and extending stub data by the matching method with those obtained by an unweighted, least-squares fit of the weighted average retirement ratios. Since the stub data was obtained by truncating complete, original life tables, the mortality behavior predicted by the two methods could be compared with the mortality behavior which actually occurred. Cowles concluded (4, p. 112):

Under the conditions adopted, i.e., the stipulations for the analysis of the retirement data, the standard assumed, and the comparison bases used, no consistent superiority was enjoyed by either the Iowa type curve method or the use of orthogonal polynomials in estimating mortality dispersion.

Scigliano (26) fitted the Weibull hazard function to the retirement ratios of the stub data developed by Cowles. The form of the Weibull hazard function used was

$$r(t) = \alpha \lambda t^{\alpha-1}$$

$r(t)$ = retirement ratio

t = time

α = shape parameter

λ = scale parameter

He used the ". . . Gauss-Newton iteration scheme for non-linear regression analysis . . ." to estimate α and λ (26, p. 21).

The stub curves were also fitted by matching (he used the results of Cowles' study) and by an unweighted least-squares fit of the weighted average retirement ratios. The matching method and the polynomial retirement ratio method appeared to yield somewhat better estimates than the Weibull hazard function (26, p. 99). However:

. . . in most of the those cases where the present methods were superior the computation method or data caused the error (26, p. 57).

Either the Gompertz or the Gompertz-Makeham equation can be fitted to the original survivor curve. The Gompertz equation is (19, p. 112)

$$L_x = kg^{c^x}$$

and the Gompertz-Makeham equation is (19, p. 113)

$$L_x = ks^x g^{c^x}$$

where

$$L_x = \text{percent surviving at age } x$$

k, s, g, c = constants to be determined from the data

The Gompertz equation expresses the "force" of retirement as an increased inability of the property to "withstand" retirement as the age of the property increases. The Gompertz-Makeham equation includes the above "force" and, also, a constant, chance "force" of retirement unrelated to age.

Nichols (23) investigated moments of the frequency distribution as means of estimating average service life and dispersion pattern. The procedure developed for estimating these mortality characteristics involved both a mathematical model and a matching process.

Two moment ratios were utilized

$$M_1 = \frac{\text{second moment of the frequency curve about the mean}}{(\text{mean})^2}$$

$$M_2 = \frac{\text{third moment of the frequency curve about the mean}}{(\text{mean})^3}$$

The standard moment ratios were calculated from the complete, standard Iowa type curves and the standard Iowa type curves stubbed at various points. Plots were made of M_1 versus M_2 at each of the percent surviving points, M_1 versus percent surviving, M_2 versus percent surviving, and M_1/M_2 versus percent surviving.

The test data used were those obtained by Cowles (4) stubbed at two different levels of percent surviving. The $M_1 - M_2$ percent surviving plot was the primary classifying plot. Where the $M_1 - M_2$ percent surviving plot did not yield a clear indication of a particular type of curve, the other plots (mentioned above) were used as supplementary guides. Service life multipliers were developed from which an estimate of the probable average service life could be obtained.

Although some of the test curves could not be classified at all, Nichols concluded that the moment ratio method appears to be a valid method of life analysis (23, p. 88).

Krane (18) developed a procedure for fitting a polynomial to the time integral of the retirement ratios and thus obtaining the life table by graduating the negative exponential function.

$$y(x) = e^{-g(x)}$$

$y(x)$ = proportion surviving at age x

$g(x)$ = time integral of the retirement ratio function

$$= \int_0^x r(t)dt$$

For large samples it is found that the covariance structure for the polynomial regression of $y(t)$ ($g(x)$ in the above notation) on t may be obtained from the multinomial distribution when the data are grouped. Thus the method of weighted least squares may be employed in fitting $y(t)$ (i.e., $g(x)$). "Censored" data in no way vitiate the method (18, p. 161).

Krane applied his procedure to one set of data from Cowles' study (4). The results were encouraging.

Henderson (13) fitted the cumulative distribution form of the Weibull function to the data of Cowles (4). The form of the Weibull cumulative distribution function utilized was (13, p. 38)

$$F(x) = 100(1 - \exp[-(L/\exp \mu)^{1/B}])$$

$\exp = e$

L = life of property

μ = position of the mode

B = scale parameter

The results of the Weibull fits of the data were compared to the results of fitting Iowa curves to the data by the matching method. The Weibull distribution yielded a better fit for symmetrical type data (13, p. 47). The Iowa curves fit data with the mode to the right of the mean better than the Weibull distribution (13, p. 47). No significant difference was found (1) in the ability of the two methods to fit data with the mode to the left of the mean (13, p. 47) and (2) in the general ability of the two methods to describe industrial property mortality experience (13, p. 57).

A smoothed life table may also be obtained by fitting a mathematical function to the survivor ratios.

$$s_{\cdot k} = \text{survivor ratio for the } k^{\text{th}} \text{ age interval}$$

$$= \frac{s_{\cdot k+1}}{s_{\cdot k}}$$

The commonly used mathematical function is a polynomial. The amount surviving at each age is calculated by the successive multiplications of the amount surviving at the beginning of the age interval by the interpolated (or extrapolated) survivor ratio for the age interval.

No one method of obtaining a smoothed, complete life table seems to yield the best results in all situations. Both the matching method (using Iowa type curves) and the retirement ratio method (using polynomials) are frequently used and yield satisfactory results in many situations.

INVESTIGATION

The vertical distribution of retirement ratios at an age interval was investigated empirically by simulation. An estimate of the form of the vertical distribution was obtained by comparing the plot of the cumulative distribution of the simulated retirement ratios with a plot of the cumulative distribution of a standard distribution function. Iowa type curves were used to provide the underlying, horizontal retirement dispersion patterns.

All calculations were performed on an IBM System/360 Model 50 digital computer at the Computation Center of Iowa State University of Science and Technology, Ames, Iowa.

Simulation of Retirement Ratios

A particular Iowa type curve and average service life specifies a probability distribution, and hence, a cumulative distribution, of ages of units at retirement. The points on the cumulative distribution were converted into integers representing a cumulative distribution of frequencies by multiplying each point by a common, appropriate multiple of ten. Then the values of the cumulative frequencies were divided into blocks of numbers representing age intervals.

The age interval during which a unit of a vintage group is retired was simulated by drawing a random number from a uniform distribution and finding that block of values (of the cumulative distribution of frequencies) which contained the random number. Additional random numbers from a uniform distribution were drawn and processed in a similar fashion to simulate the retirement of all units in the vintage group. Since the

retirement ratio for the k^{th} age interval is

$$r_k = \frac{\text{number of units retired during the } k^{\text{th}} \text{ age interval}}{\text{number of units surviving at the beginning of the } k^{\text{th}} \text{ age interval}}$$

the retirement ratio for each age interval could be calculated. In this manner, a set of retirement ratios, one for each age interval, for a vintage group was simulated.

Additional simulation runs using the same size vintage group (generally the retirement experience of a vintage group was simulated 100 times) and based on the same Iowa type curve and average service life, yield additional sets of retirement ratios, each set containing one retirement ratio for each age interval. The retirement ratios at each age interval, one from each set, represent a sample (of size equal to the number of simulation runs) from the population of retirement ratios for that age interval from a vintage group of the given size.

The property group size was specified in terms of physical units rather than dollars. The reason for using physical units is that the age at retirement of any one physical unit is independent of the age at retirement of any other physical unit. Dollars (of property) do not have this independence unless each dollar represents exactly one physical unit.

The retirement ratios at each age interval were arranged in ascending order. A cumulative count of the ordered retirement ratios at an age interval yielded the cumulative distribution of retirement ratios at that age interval. The cumulative counts at each age interval were converted to cumulative percents to facilitate comparison of the simulated cumulative distributions with the cumulative distribution of a standard distribution

function. A flow chart of the computer program developed to simulate the vertical distributions of retirement ratios is shown in Appendix A.

Simulation runs based on different parent populations (i.e., Iowa type curve, average service life, and property group size) yielded additional sets of vertical distributions of retirement ratios.

Normal Approximation

The simulated, vertical distributions of retirement ratios were plotted on both rectangular co-ordinate graph paper and normal probability paper. The cumulative distribution points for each age interval plotted on normal probability paper lie closely about a straight line, except the points for the early and late age intervals. Hence, the normal distribution appeared to be a likely candidate for representing the vertical distribution of retirement ratios at an age interval.

For a specified Iowa type curve, average service life and property group size, let

$$L_{ijk} = 1 \text{ if the } j^{\text{th}} \text{ unit of the } i^{\text{th}} \text{ simulation run is retired} \\ \text{during the } k^{\text{th}} \text{ age interval} \\ = 0 \text{ otherwise}$$

$$M_{ijk} = 1 \text{ if the } j^{\text{th}} \text{ unit of the } i^{\text{th}} \text{ simulation run is retired} \\ \text{after the } k^{\text{th}} \text{ age interval} \\ = 0 \text{ otherwise}$$

$$Z_{ijk} = 1 \text{ if the } j^{\text{th}} \text{ unit of the } i^{\text{th}} \text{ simulation run is retired} \\ \text{before the } k^{\text{th}} \text{ age interval} \\ = 0 \text{ otherwise}$$

i = simulation run index number

$= 1, 2, \dots, I$

j = property unit index number

$= 1, 2, \dots, J$

k = age interval index number

$= 1, 2, \dots, K$

r_{ik} = retirement ratio for the k^{th} age interval of the i^{th} simulation run

Then

$$r_{ik} = \frac{\sum_{j=1}^J L_{ijk}}{\sum_{j=1}^J L_{ijk} + \sum_{j=1}^J M_{ijk}}$$

$$= \frac{L_{i \cdot k}}{L_{i \cdot k} + M_{i \cdot k}}$$

where

$$L_{i \cdot k} = \sum_{j=1}^J L_{ijk}$$

$$M_{i \cdot k} = \sum_{j=1}^J M_{ijk}$$

If

$$\Pr(L_{ijk} = 1) = C_k; j = 1, 2, \dots, J; i = 1, 2, \dots, I$$

$$\Pr(M_{ijk} = 1) = C'_k; j = 1, 2, \dots, J; i = 1, 2, \dots, I$$

$$\Pr(Z_{ijk} = 1) = C''_k; j = 1, 2, \dots, J; i = 1, 2, \dots, I$$

the $L_{i \cdot k}$, $M_{i \cdot k}$, and $Z_{i \cdot k}$ are each binomially distributed and collectively they form a multinomial distribution; C_k , C'_k , and C''_k are the corresponding probabilities of the multinomial distribution.

The cumulative distribution of r_{ik} , for some specified k , can be obtained by calculating the probabilities

$$\Pr\left(\frac{L_{i \cdot k}}{L_{i \cdot k} + M_{i \cdot k}} \leq T\right)$$

as the dummy variable T is varied from zero to one, the range of a retirement ratio. A more useful form of the above expression is

$$\begin{aligned} \Pr\left(\frac{L_{i \cdot k}}{L_{i \cdot k} + M_{i \cdot k}} \leq T\right) &= \Pr(L_{i \cdot k} \leq T L_{i \cdot k} + T M_{i \cdot k}) \\ &= \Pr[(1 - T)L_{i \cdot k} - T M_{i \cdot k} \leq 0] \end{aligned}$$

The mean of the expression

$$[(1 - T)L_{i \cdot k} - T M_{i \cdot k}]$$

is the expected value of the expression. Therefore

$$\begin{aligned} E[(1 - T)L_{i \cdot k} - T M_{i \cdot k}] &= E[(1 - T)L_{i \cdot k}] - E[T M_{i \cdot k}] \\ &= (1 - T)E(L_{i \cdot k}) - T E(M_{i \cdot k}) \end{aligned}$$

since (1, p. 32):

1. The expected value of a sum (or difference) of two variates or functions is the sum (or difference) of the expected values of the separate parts.
2. The expected value of a constant times a variable is the constant times the expected value of the variable.

Both $L_{i \cdot k}$ and $M_{i \cdot k}$ are, individually, binomially distributed. The mean of a variable which is binomially distributed is usually expressed as (1, p. 35)

$$\mu = np$$

where

n = number of independent trials and is analogous to J

p = probability of success on any one trial and is analogous to

$$C_k \text{ and } C'_k$$

Therefore

$$E(L_{i \cdot k}) = J C_k$$

$$E(M_{i \cdot k}) = J C'_k$$

where, as mentioned above

$$C_k = \Pr(L_{ijk} = 1)$$

$$C'_k = \Pr(M_{ijk} = 1)$$

Then

$$E[(1 - T)L_{i \cdot k} - T M_{i \cdot k}] = (1 - T) J C_k - T J C'_k$$

The variance of the expression

$$[(1 - T)L_{i \cdot k} - T M_{i \cdot k}]$$

is the expected value of the square of a similar expression but where the mean of each variable is subtracted from the variable.

$$\begin{aligned} \text{var}[(1 - T)L_{i \cdot k} - T M_{i \cdot k}] \\ = E\{[(1 - T)(L_{i \cdot k} - \mu_k) - T(M_{i \cdot k} - \mu'_k)]^2\} \end{aligned}$$

where

$$\mu_k = E(L_{i \cdot k})$$

$$\mu'_k = E(M_{i \cdot k})$$

Then

$$\begin{aligned} \text{var}[(1 - T)L_{i \cdot k} - T M_{i \cdot k}] \\ = E[(1 - T)^2(L_{i \cdot k} - \mu_k)^2 + T^2(M_{i \cdot k} - \mu'_k)^2 \\ - 2(T)(1 - T)(L_{i \cdot k} - \mu_k)(M_{i \cdot k} - \mu'_k)] \\ = (1 - T)^2 E[(L_{i \cdot k} - \mu_k)^2] + T^2 E[(M_{i \cdot k} - \mu'_k)^2] \\ - 2(1 - T)T E[(L_{i \cdot k} - \mu_k)(M_{i \cdot k} - \mu'_k)] \end{aligned}$$

The variance of a binomial is (1, p. 35)

$$\begin{aligned} \mu_2 &= E(y - \mu)^2 \\ &= npq \end{aligned}$$

where

n, p = as previously defined

$$q = 1 - p$$

The covariance of two of the variates of a multinomial is (1, p. 50, 54)

$$\begin{aligned}\mu_{11} &= E[(y_1 - \mu_{10})(y_2 - \mu_{01})] \\ &= -np_{10}p_{01}\end{aligned}$$

where

p_{10} = probability of the event y_1 occurring and is analogous to

$$c_k$$

p_{01} = probability of the event y_2 occurring and is analogous to

$$c'_k$$

Therefore

$$\begin{aligned}E[(L_{i.k} - \mu_k)^2] &= J c_k(1 - c_k) \\ E[(M_{i.k} - \mu'_k)^2] &= J c'_k(1 - c'_k) \\ E[(L_{i.k} - \mu_k)(M_{i.k} - \mu'_k)] &= -J c_k c'_k\end{aligned}$$

Then

$$\begin{aligned}\text{var}[(1 - T)L_{i.k} - T M_{i.k}] &= (1 - T)^2 J c_k(1 - c_k) + T^2 c'_k(1 - c'_k) \\ &\quad + 2T(1 - T) J c_k c'_k\end{aligned}$$

Since the multinomial distribution is well approximated by the normal distribution, at least for non-extreme parametric values, the linear combination of multinomial counts, $(1 - T) L_{i.k} - T M_{i.k}$, will also be approximated by the normal distribution; it is to be expected that this linear combination will be more nearly normal than the ratio $L_{i.k}/(L_{i.k} - M_{i.k})$. The distribution of the ratio is therefore approximated using the approximate normality of the linear combination:

$$\begin{aligned}
& \Pr\left(\frac{L_{i \cdot k}}{L_{i \cdot k} + M_{i \cdot k}} \leq T\right) \\
&= \Pr[(1 - T)L_{i \cdot k} - T M_{i \cdot k} \leq 0] \\
&\doteq \Pr[N(\mu, \sigma) \leq 0] \\
&\doteq \Pr\{N[(1 - T) J C_k - T J C'_k, \{(1 - T)^2 J C_k(1 - C_k) \\
&\quad + T^2 J C'_k(1 - C'_k) + 2T(1 - T) J C_k C'_k\}^{1/2}] \leq 0\}
\end{aligned}$$

A digital computer program for computing the points of the normal cumulative distribution, based on an approximation by Hastings (12, p. 168), was obtained from the Iowa State University Statistical Laboratory - Numerical Analysis and Programming Section, Ames, Iowa. Utilizing this program and the theoretical values of C_k and C'_k (based on the Iowa type curve and average service life) and the vintage group size, a pseudo-normal, cumulative distribution of retirement ratios at each age interval was calculated (see Appendix B).

RESULTS OF THE INVESTIGATION

The first program to calculate the points of the pseudo-normal, cumulative distributions utilized theoretical C_k and C'_k values which were dependent upon the chosen vintage group size. Plots of the pseudo-normal, cumulative distributions on rectangular co-ordinate graph paper did not satisfactorily match with plots of the simulated cumulative distributions; the general shapes of the pseudo-normal, cumulative distribution plots were appropriate but the (horizontal) locations were not.

The C_k and C'_k values were calculated in the following manner in the first program. The "theoretical" number of units surviving at each age from a vintage group of the chosen size (for the chosen Iowa type curve and average service life) were calculated from a table¹ of the theoretical percent surviving; all values were rounded to the nearest whole unit. Then

$$\begin{aligned} C_k &= \Pr(L_{ijk} = 1) \\ &= \frac{\text{number of units retired during age interval } k}{\text{vintage group size}} \end{aligned}$$

$$\begin{aligned} C'_k &= \Pr(M_{ijk} = 1) \\ &= \frac{\text{number of units retired after age interval } k}{\text{vintage group size}} \end{aligned}$$

Thus, these estimates of the theoretical C_k and C'_k were dependent upon vintage group size, curve type and average service life rather than just the curve type and average service life.

¹Scigliano, J. Michael, Graduate Assistant in the Department of Industrial Engineering, Iowa State University of Science and Technology, Ames, Iowa. Tables of the theoretical percent surviving, to six decimal places, at 1% intervals of an average service life of 100 years for the Iowa type curves based on the original data of Robley Winfrey. Private communication. 1965.

A minor modification of the first program, calculating C_k and C'_k on the basis of a "parent population" of 100,000,000 units rather than the chosen vintage group size, yielded more successful results. Although the theoretical "parent population" is not limited to 100,000,000 units, the amount of error introduced in estimating C_k and C'_k by the use of such a large number of units is less than when a small vintage group size is used (and is probably negligible).

Retirement ratios of the forms (1) zero divided by zero and (2) zero divided by a positive number, occurred in the simulations of the retirement ratios and were assigned a value of zero. The first form arose whenever the age interval during which the oldest unit of a particular sample (a particular simulation run of the retirement experience of the vintage group) was retired was, say, k and the age interval during which the oldest unit of some other sample was retired was, say, $k + 1$. The retirement ratio from the first-mentioned sample for the $(k + 1)^{st}$ age interval would, then, have to be of the form zero divided by zero. The second form occurred whenever one or more units were surviving at the beginning of an age interval and no units were retired during that age interval.

Retirement ratios of the form zero divided by zero could not occur in the pseudo-normal program because the computer was programmed to stop whenever the number of units surviving at the beginning of an age interval was zero. A retirement ratio of the form zero divided by a positive number could possibly occur in the pseudo-normal program only in the early age intervals when the theoretical values of C_k and C'_k were

$$C_k = 0$$

$$C'_k = 1$$

The pseudo-normal program assigned a value of one to all of the cumulative probabilities for that age interval.

$$\Pr[N(\mu, \sigma) \leq 0] = 1 \quad 0 \leq T \leq 1$$

The pseudo-normal program computed the values of the cumulative distribution for values of T between zero and one in increments of 0.01. A large sample size coupled with a small C_k value, a large C'_k value, and a delta T of 0.01 resulted in very few cumulative distribution values other than zero or one, if any. A larger number of non-zero, non-one cumulative distribution values could have been obtained by incrementing T by 0.001 or an even smaller amount.

The vertical distributions of retirement ratios at each age interval were simulated (and the corresponding pseudo-normal, cumulative distributions calculated) for only a small number of the possible combinations of curve type, average service life, and sample size. Representative plots, on normal probability paper, of the simulated cumulative distributions of retirement ratios and of the pseudo-normal, cumulative distributions of retirement ratios are shown in Figures 9a, 9b, 10a, and 10b for an Iowa $L_3 - 10$ and in Figures 11a, 11b, 12a, and 12b for an Iowa $R_1 - 25$. For both the $L_3 - 10$ and $R_1 - 25$ simulations, the sample size (vintage group size) was 100 and the retirement experience of a sample was simulated 100 times.

A few general comments can be made from a visual inspection of the probability plots (including those not shown herein):

Figure 9a. Simulated cumulative distributions of retirement ratios for $L_3 - 10$

- - age interval 1.5 - 2.5 years
- - age interval 4.5 - 5.5 years
- △ - age interval 7.5 - 8.5 years
- ◇ - age interval 10.5 - 11.5 years

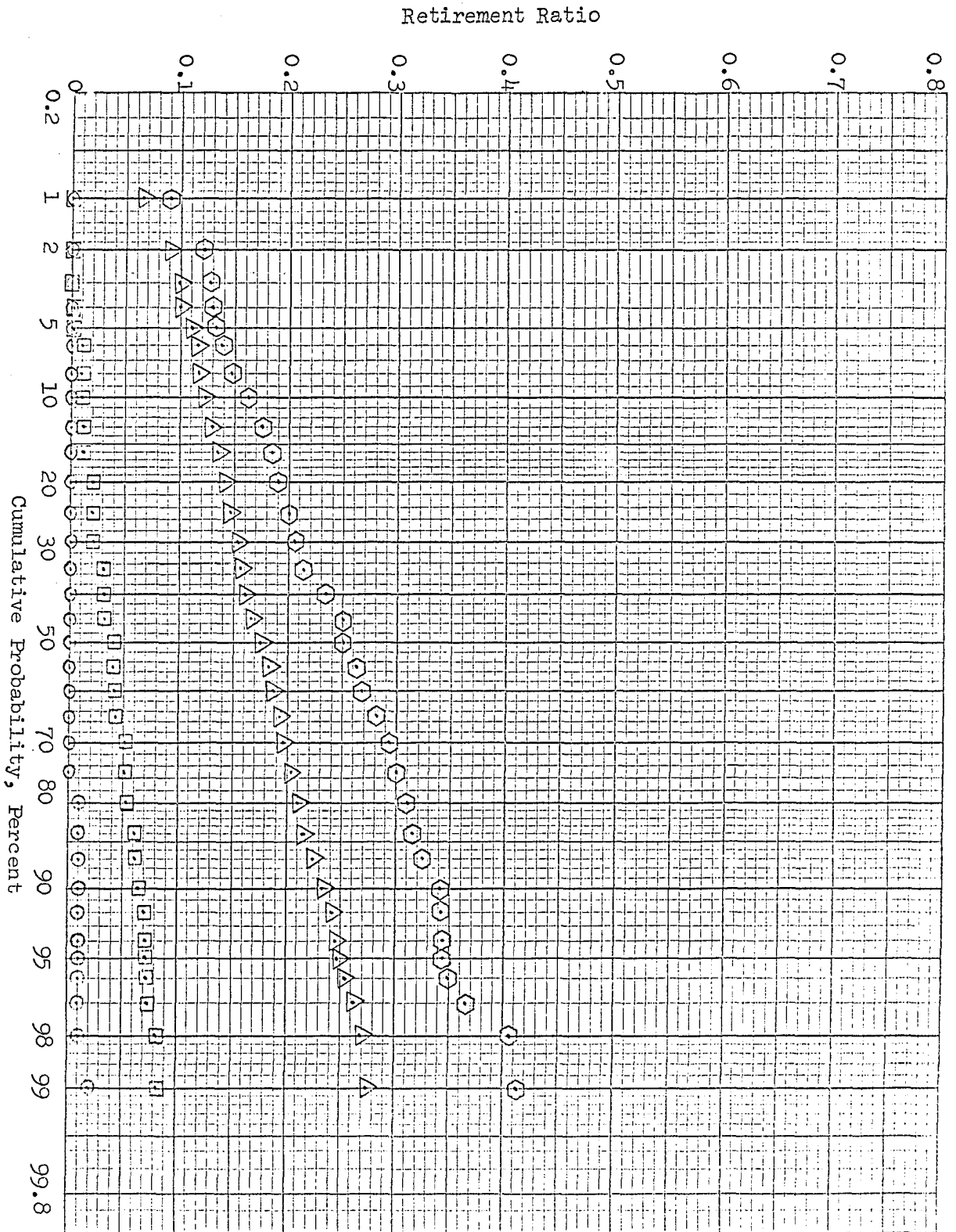


Figure 9b. Simulated cumulative distributions of retirement ratios for $L_3 - 10$

○ - age interval 13.5 - 14.5 years

□ - age interval 16.5 - 17.5 years

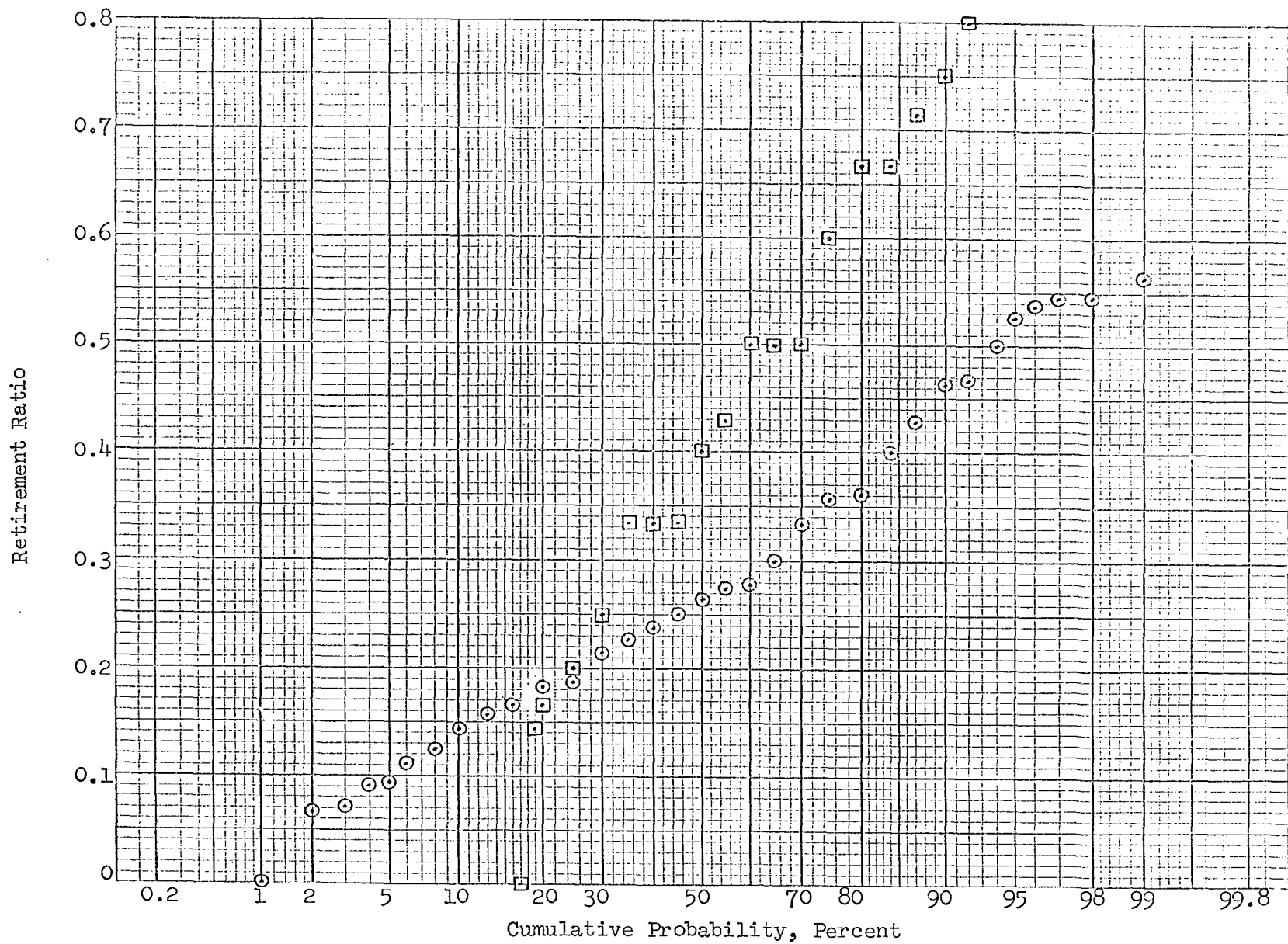


Figure 10a. Pseudo-normal cumulative distributions of T for $L_3 - 10$

○ - age interval 1.5 - 2.5 years

□ - age interval 4.5 - 5.5 years

△ - age interval 7.5 - 8.5 years

◇ - age interval 10.5 - 11.5 years

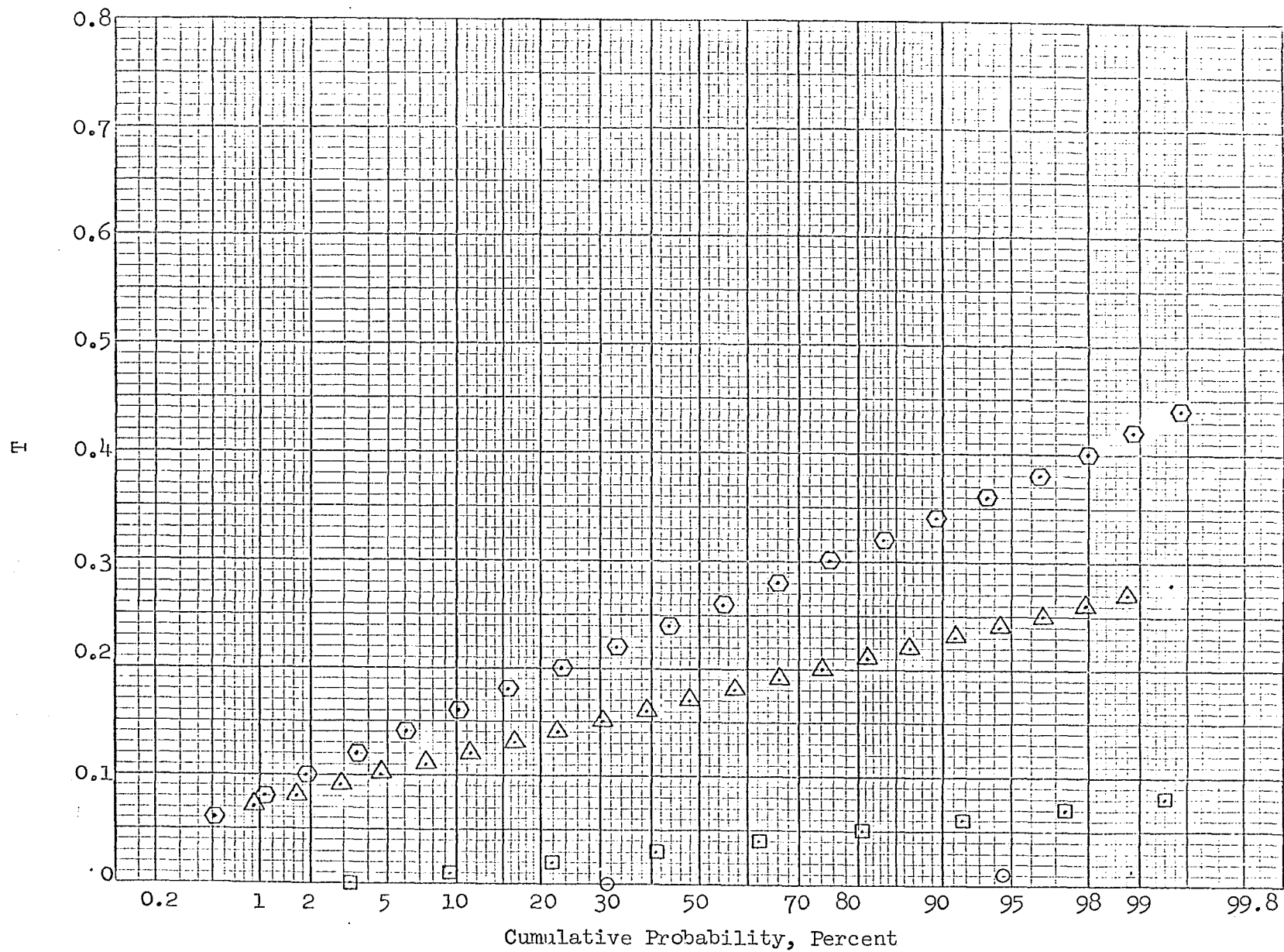


Figure 10b. Pseudo-normal cumulative distribution of T for $L_3 - 10$

○ - age interval 13.5 - 14.5 years

□ - age interval 16.5 - 17.5 years

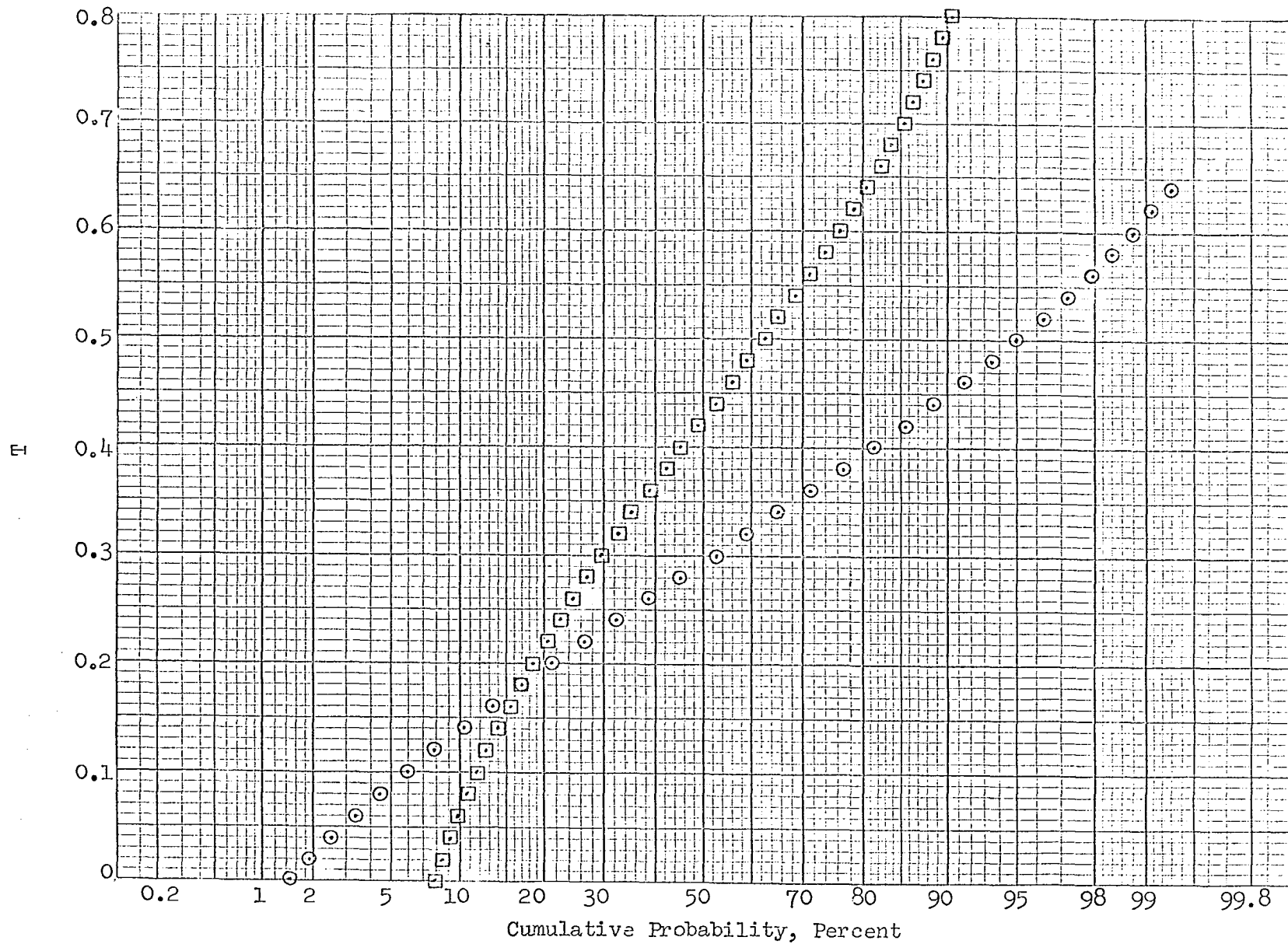


Figure 11a. Simulated cumulative distributions of retirement ratios for $R_1 - 25$

- - age interval 4.5 - 5.5 years
- - age interval 12.5 - 13.5 years
- △ - age interval 18.5 - 19.5 years
- ⊙ - age interval 24.5 - 25.5 years

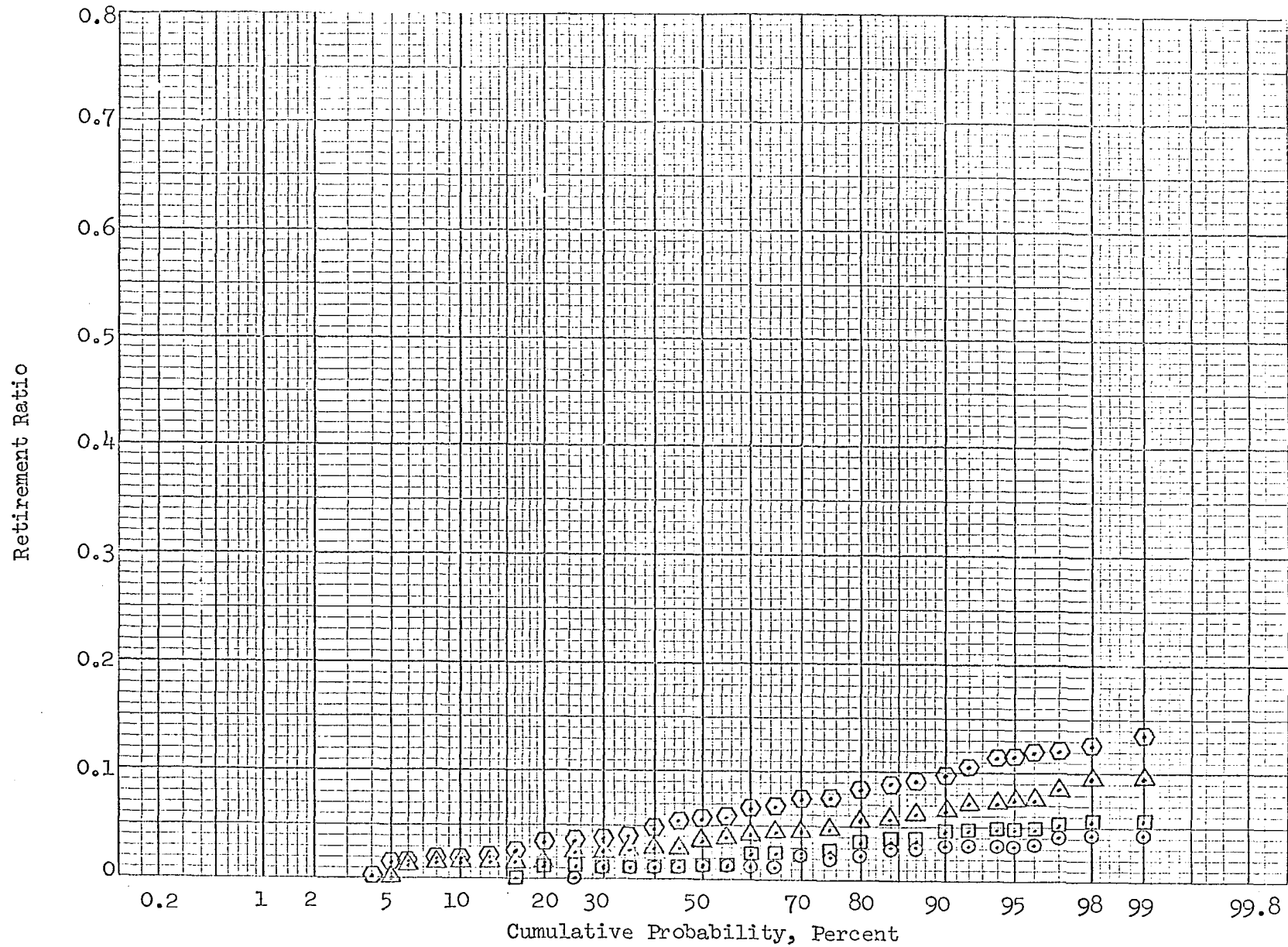


Figure 11b. Simulated cumulative distributions of retirement ratios for $R_1 - 25$

○ - age interval 30.5 - 31.5 years

□ - age interval 36.5 - 37.5 years

◇ - age interval 44.5 - 45.5 years

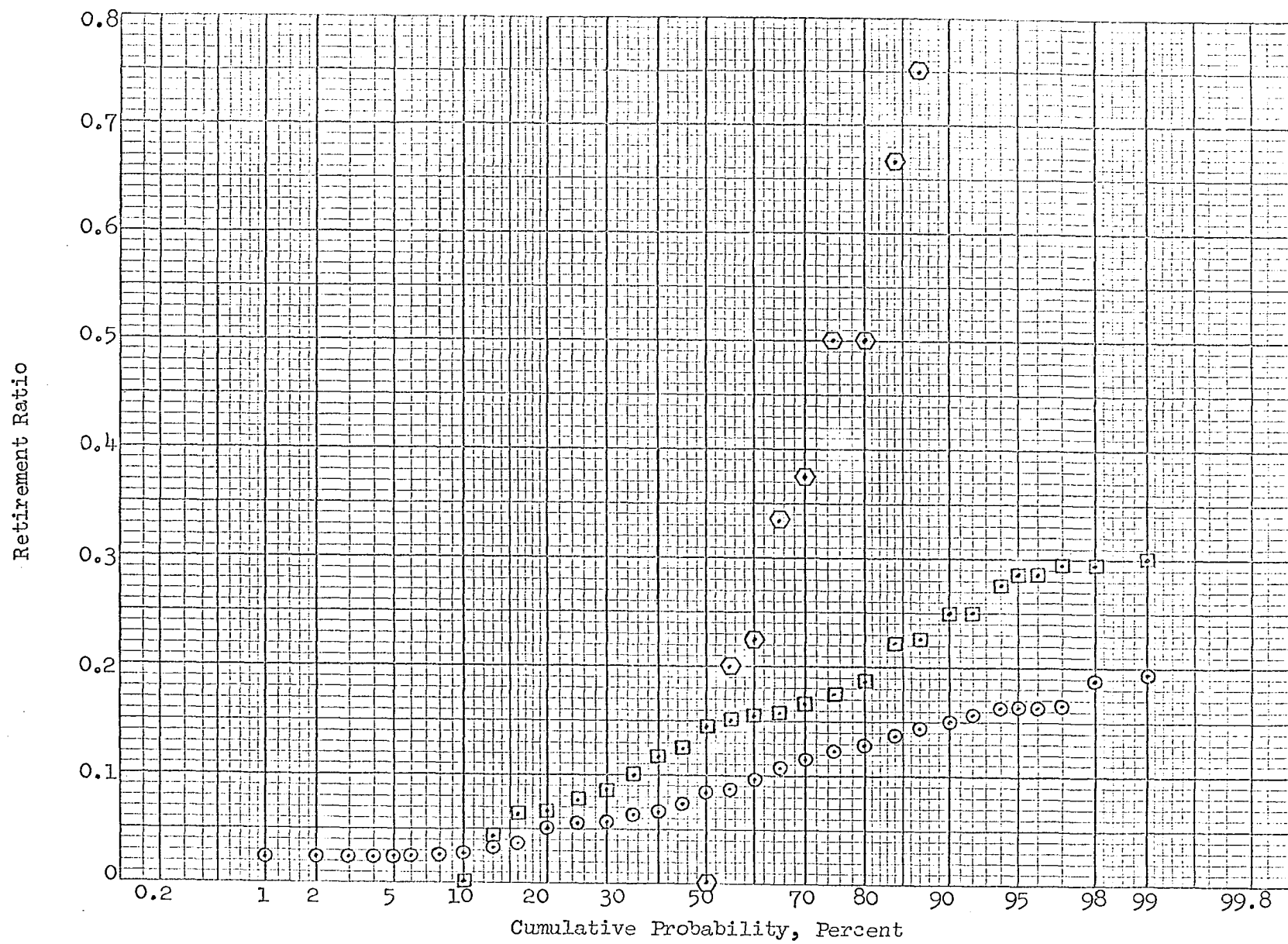


Figure 12a. Pseudo-normal cumulative distributions of T for $R_1 - 25$

○ - age interval 4.5 - 5.5 years

□ - age interval 12.5 - 13.5 years

△ - age interval 18.5 - 19.5 years

○ - age interval 24.5 - 25.5 years

E-I

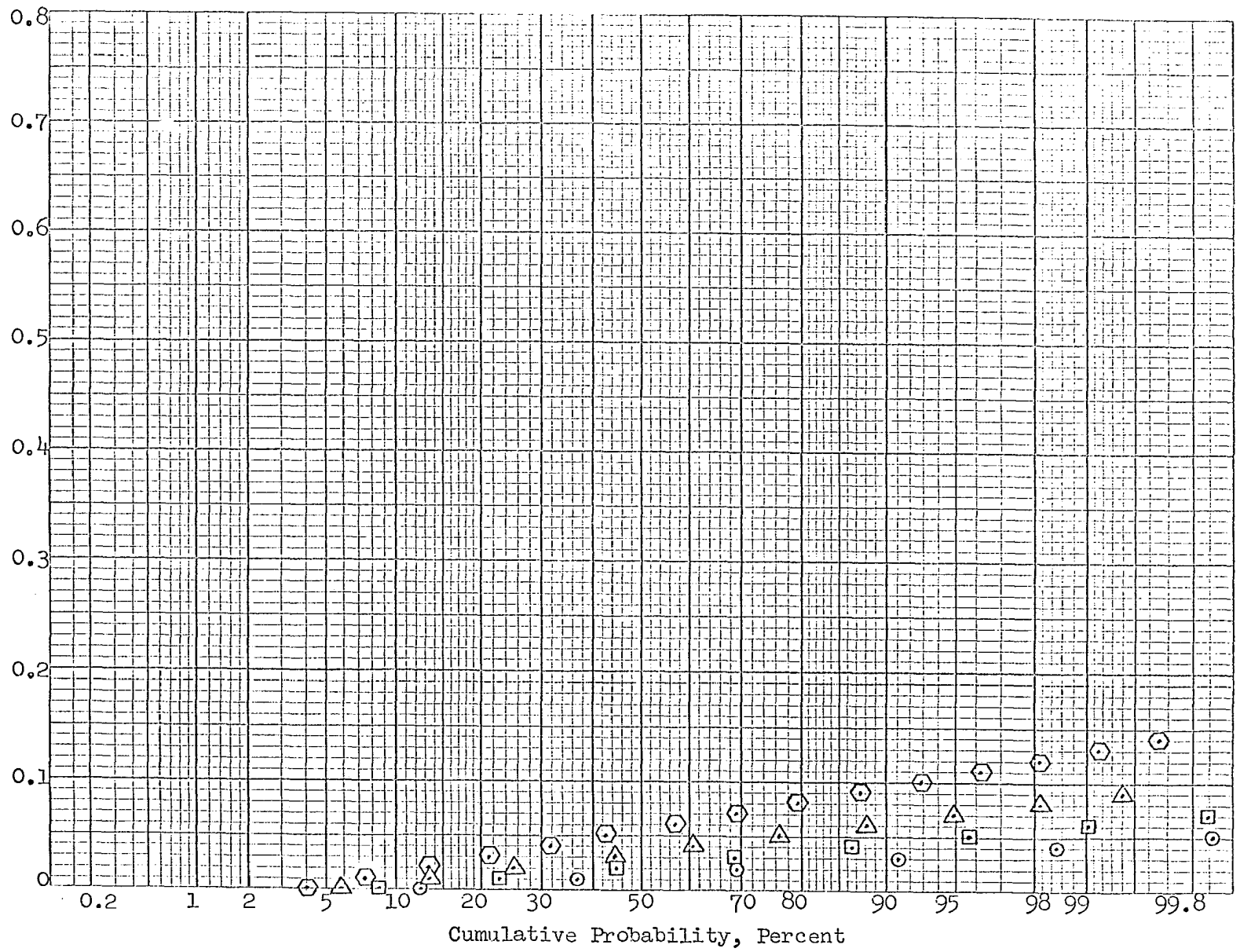
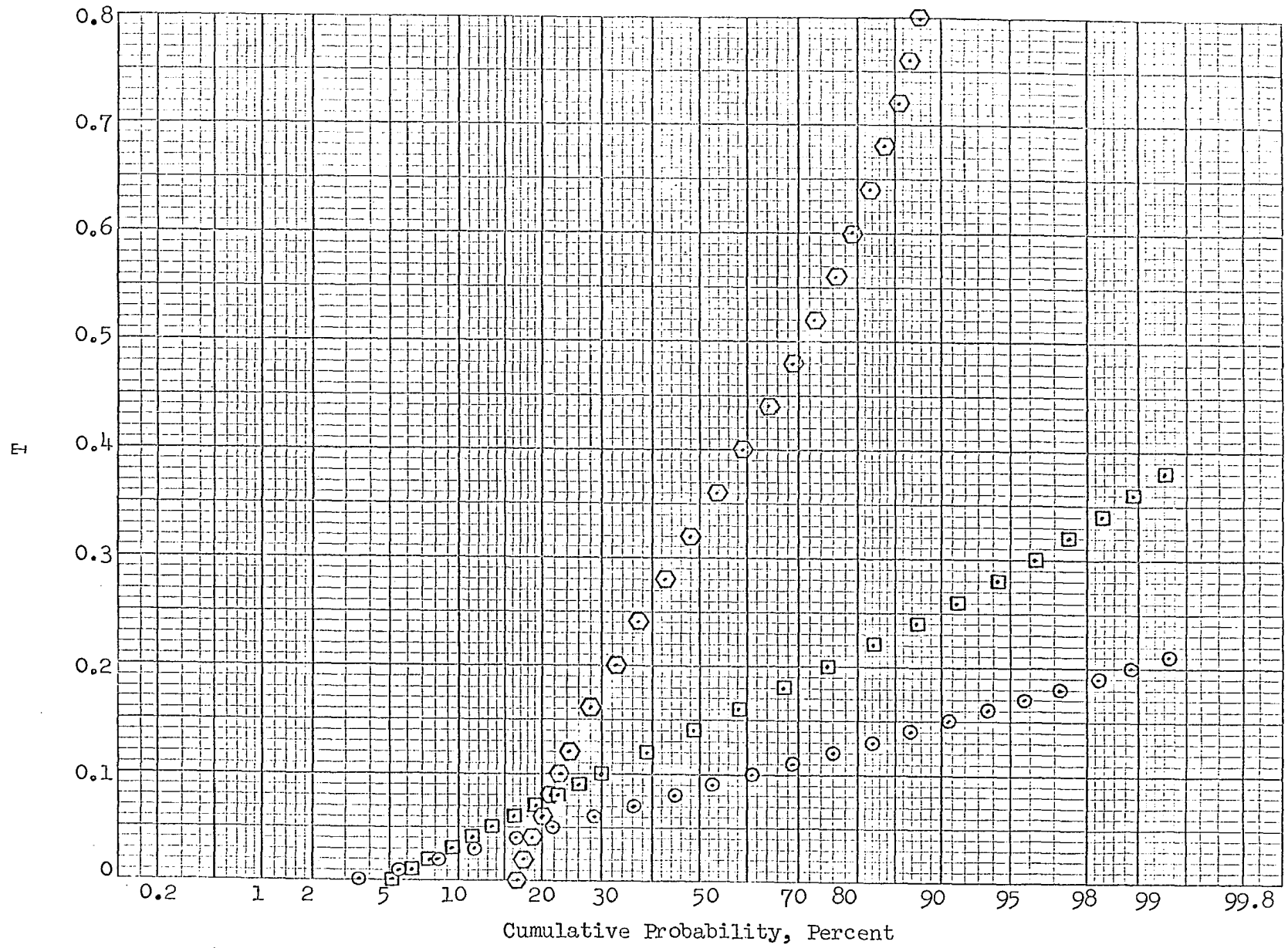


Figure 12b. Pseudo-normal cumulative distributions of T for $R_1 - 25$

○ - age interval 30.5 - 31.5 years

□ - age interval 36.5 - 37.5 years

◇ - age interval 44.5 - 45.5 years



1. The slope of a straight line drawn through the points for an age interval reflects the standard deviation of the sample and increases as the age interval index number increases (for a given Iowa type curve, average service life, and sample size).
2. The slope of a straight line drawn through the points for a given age interval decreases as the sample size increases (for a given Iowa type curve, and average service life).

The three most important results of the investigation were:

1. The variance of the vertical distribution of retirement ratios does not remain constant from age interval to age interval.
2. The points of the simulated cumulative distribution of retirement ratios at each age interval, plotted on normal probability paper, lie along or nearly along a straight line, except for the early and late age intervals.
3. Plots of the pseudo-normal, cumulative distributions of retirement ratios at each age interval match satisfactorily, visually, with the plots of the simulated cumulative distributions of retirement ratios at each age interval, except for late age intervals.

The first result indicates that the assumption of homoscedasticity (constant variance), a necessary condition if the unweighted least-squares method is to yield linear unbiased estimators which have minimum variance amongst the class of linear unbiased estimators of the polynomial coefficients, is invalid. The second result indicates that the vertical distribution of retirement ratios at an age interval is approximately a normal distribution. The third result permits the computation of an estimate of the variance of a retirement ratio, by means of the pseudo-normal

approximation, when only the vintage group size and estimates of C_k and C'_k are known.

A PROCEDURE FOR FITTING A POLYNOMIAL TO RETIREMENT RATIOS

A procedure for fitting a polynomial to the retirement ratios calculated from original data is developed in this section. The procedure is based on the least-squares principle.

Assumptions

A number of the assumptions necessary to make an actuarial life analysis have already been presented in previous sections of this dissertation. The procedure herein developed for fitting a polynomial to the retirement ratios is dependent on the above mentioned assumptions plus a few additional assumptions, all of which are listed below:

1. Basic assumptions of life analysis:
 - a. The mortality behavior of the property follows some "law of mortality" expressible as a function of time.
 - b. The past mortality behavior of a property is indicative of the expected future mortality behavior of the property.
2. Assumptions concerning the property data used:
 - a. The data available is from the historical records of the same property or a similar property.
 - b. A life analysis based on physical units is meaningful.
 - c. Sufficient data in a usable form are available to make an actuarial life analysis study.
 - d. The data set selected for analysis is composed of homogeneous units (both within and between vintage groups) which follow the same "law of mortality".

3. Assumptions in identifying the law of mortality:
 - a. The "law of mortality" is better represented by a smooth curve fitted to raw data points than by unsmoothed, raw data points.
 - b. The "law of mortality" can be adequately represented by a polynomial expressing the relationship between retirement ratios and age intervals.
4. Assumptions of the method of fitting a polynomial to the retirement ratios:
 - a. The age interval of retirement of those units already retired is determined without error.
 - b. The regression of retirement ratios on age intervals is linear in the polynomial coefficients.
 - c. The retirement ratios calculated from the data of a given vintage group are independent of each other (i.e., the assumption from (16, p. 383): "The deviations $y_j - E(y|x_j)$ are mutually independent").
 - d. The vintage groups contributing mortality experience to the data set are independent, random samples (but not necessarily of the same size) from the same parent population of physical units.
 - e. Each retirement ratio is a random sample from an approximately normally distributed parent population of retirement ratios with parameters (1) curve type, (2) average service life, (3) sample size, and (4) age interval, and the distribution

of the parent population can be approximated by a pseudo-normal, cumulative distribution (the assumption that the distribution is normal is not needed to develop the procedure but is needed in testing the significance of the n^{th} degree term of the retirement ratio polynomial).

- f. The expected value of a vertical distribution of retirement ratios is a constant for a given curve type, average service life, and age interval, regardless of sample size.

Development of Procedure

Under the assumptions of the method of fitting a polynomial to the retirement ratios (above), the principle of least-squares yields certain estimators of the polynomial coefficients. The properties of these estimators are dependent upon the assumptions that can reasonably be made about the r_{ik} and the variances of the r_{ik} . Let

μ_k = population retirement ratio for the k^{th} age interval

r_{ik} = sample retirement ratio for the k^{th} age interval from the i^{th} vintage group

σ_{ik}^2 = variance of the population of retirement ratios for the k^{th} age interval for samples of the size of vintage group i

$\hat{\sigma}_{ik}^2$ = sample estimate of σ_{ik}^2

$r_{\cdot k}$ = a weighted average (over all i) of the r_{ik}

Because of computational considerations, the r_{ik} for each age interval must be combined into a single $r_{\cdot k}$. Since

$$E(r_{ik}) = \mu_k; i = 1, \dots, I$$

$$\text{var}(r_{ik}) = \sigma_{ik}^2$$

$r_{.k}$, the single-valued, least-squares estimator, is (11, p. 12)

$$\begin{aligned} r_{.k} &= \frac{\sum_{ik}^I w_{ik} r_{ik}}{\sum_{ik} w_{ik}} \\ &= \frac{\sum_{ik}^I (\sigma^2/\sigma_{ik}^2) r_{ik}}{\sum_{ik}^I \sigma^2/\sigma_{ik}^2} \\ &= \frac{\sum_{ik}^I (1/\sigma_{ik}^2) r_{ik}}{\sum_{ik}^I 1/\sigma_{ik}^2} \end{aligned}$$

where

σ = a constant of proportionality which may be used to adjust the magnitude of the relative weights

Substituting the sample estimate of σ_{ik}^2 , $\hat{\sigma}_{ik}^2$, yields

$$r_{.k} = \frac{\sum_{ik}^I (1/\hat{\sigma}_{ik}^2) r_{ik}}{\sum_{ik}^I 1/\hat{\sigma}_{ik}^2}$$

Then, the least-squares expression to minimize is (11, p. 88)

$$\text{Min}_{a,b,c,\text{etc.}} \quad 1/\sigma^2 \sum_{ik}^K w_{.k} [r_{.k} - (a + bx_k + cx_k^2 + \dots)]^2$$

where

$$\begin{aligned} w_{.k} &= 1/\text{var}(r_{.k}) \\ &= \sum_{ik}^I 1/\sigma_{ik}^2 \end{aligned}$$

Differentiation of the least-squares expression yields the normal equations, of the form

$$\sum_{k=1}^K w_{.k} [r_{.k} - (a + bx_k + cx_k^2 + \dots)] = 0$$

$$\sum_{k=1}^K w_{.k} [r_{.k} - (a + bx_k + cx_k^2 + \dots)]x = 0$$

$$\sum_{k=1}^K w_{.k} [r_{.k} - (a + bx_k + cx_k^2 + \dots)]x_k^2 = 0$$

etc.

Replacing σ_{ik}^2 by $\tilde{\sigma}_{ik}^2$, the sample estimate of σ_{ik}^2 , the first normal equation becomes

$$\sum_{k=1}^K \sum_{i=1}^I 1/\tilde{\sigma}_{ik}^2 [r_{.k} - (a + bx_k + cx_k^2 + \dots)] = 0$$

or

$$\sum_{k=1}^K \sum_{i=1}^I 1/\tilde{\sigma}_{ik}^2 \left[\frac{\sum_{i=1}^I 1/\tilde{\sigma}_{ik}^2 r_{ik}}{\sum_{i=1}^I 1/\tilde{\sigma}_{ik}^2} - (a + bx_k + cx_k^2 + \dots) \right] = 0$$

or

$$\sum w_{.k} r_{.k} = a \sum w_{.k} + b \sum w_{.k} x_k + c \sum w_{.k} x_k^2 + \dots$$

and similarly for the other normal equations.

The maximum-likelihood estimators of the polynomial coefficients, under the additional assumption

$$r_{ik} \sim N(\mu_k, \sigma_{ik})$$

yields the same set of equations as the principle of least squares (see Appendix C).

A theoretical procedure for fitting a polynomial to the retirement ratios is:

1. Compute each r_{ik} .
2. Compute $\tilde{\sigma}_{ik}^2$ for each r_{ik} .
3. Fit a polynomial to the r_{ik} by the weighted least-squares method, where the weight to be given to each r_{ik} is $w_{ik} = 1/\tilde{\sigma}_{ik}^2$.

The variance, $\hat{\sigma}_{ik}^2$, obtained by the pseudo-normal, cumulative distribution is a function of both the vintage group size at age zero, J_i , and the sample retirement ratio, r_{ik} , as an estimate of the mean population retirement ratio, μ_k . If a better estimate of μ_k could be found, a better estimate of $\hat{\sigma}_{ik}^2$ could be calculated.

Assumptions 3a, 3b, and 4f indicate a way of obtaining a better estimate of μ_k to use in computing $\hat{\sigma}_{ik}^2$. Firstly, since each r_{ik} , for a given k , is assumed to be an estimate of μ_k , some average or weighted average value of the r_{ik} , say $\tilde{r}_{.k}$, should be a better estimate of μ_k . Secondly, since the "law of mortality" is assumed to be representable by a smooth curve, and in particular a polynomial function, the $\tilde{r}_{.k}$ interpolated from a polynomial function fitted to average or weighted average values of the r_{ik} should be better estimates of the μ_k than the r_{ik} .

Several methods of obtaining the $\tilde{r}_{.k}$ are available; which method is the "best" has not been established. The chosen method utilizes a preliminary approach to the over-all problem of fitting a polynomial to the retirement ratios (see Appendix D) and is as follows:

1. Compute $\tilde{r}_{.k} = \frac{L_{.k}}{L_{.k} + M_{.k}}$
2. Compute $\tilde{w}_{.k} = \frac{L_{.k} + M_{.k}}{\tilde{r}_{.k}(1 - \tilde{r}_{.k})}$
3. Fit a polynomial to the $\tilde{r}_{.k}$ by the weighted least-square method, where the weight to be given each $\tilde{r}_{.k}$ is $\tilde{w}_{.k}$.
4. Interpolate the necessary values of $\tilde{r}_{.k}$ from the polynomial of (3).

The Procedure

The over-all procedure developed for fitting a polynomial to the retirement ratios is, then,:

1. Compute $\tilde{r}_{.k} = \frac{L_{.k}}{L_{.k} + M_{.k}}$ for all k.
2. Compute $\tilde{w}_{.k} = \frac{L_{.k} + M_{.k}}{\tilde{r}_{.k}(1 - \tilde{r}_{.k})}$ for all k.
3. Fit a polynomial to the $\tilde{r}_{.k}$ by the weighted least-squares method, where the weight to give each $\tilde{r}_{.k}$ is $\tilde{w}_{.k}$.
4. Interpolate the $\tilde{r}_{.k}$ from the polynomial (3).
5. Compute $\tilde{\sigma}_{ik}^2$ for each r_{ik} from the pseudo-normal, cumulative distribution based on J_i and C_k and C'_k (from $\tilde{r}_{.k}$).
6. Fit a polynomial to the r_{ik} by the weighted least-squares method, where the weight to give each r_{ik} is $1/\tilde{\sigma}_{ik}^2$.

A procedure equivalent to step (6) is to fit a polynomial to the weighted average retirement ratio, $\hat{r}_{.k}$, at each age interval by the weighted least-squares method, where

$$\hat{r}_{.k} = \frac{\sum_I (1/\tilde{\sigma}_{ik}^2) r_{ik}}{\sum_I (1/\tilde{\sigma}_{ik}^2)}$$

$\tilde{\sigma}_{ik}^2$ = sample estimate of the variance of r_{ik} computed from the pseudo-normal, cumulative distribution

and the weight to be given each $\hat{r}_{.k}$ is

$$w_{.k} = \sum_I (1/\tilde{\sigma}_{ik}^2)$$

A general flow chart of a computer program to implement the procedure is shown in Appendix E.

Comments

An estimate of the variance of each r_{ik} , $\hat{\sigma}_{ik}^2$, can be calculated in the following manner:

1. Generate the pseudo-normal, cumulative distribution of r_{ik} based on J_i (the size of vintage group i) and C_k and C'_k computed from $\tilde{r}_{\cdot k}$.
2. Compute the area above the cumulative distribution versus T curve (area above the curve and below a horizontal line representing a cumulative probability of one).
3. Compute the area above the cumulative distribution versus T^2 curve.
4. The difference between the area computed in (3) and the square of the area computed in (2) is an estimate of $\hat{\sigma}_{ik}^2$.

The proof of (4) is as follows.

$$\begin{aligned}
 \text{var}(t) &= E(T - \mu)^2 \\
 &= E(T^2 - 2\mu T + \mu^2) \\
 &= E(T^2) - 2\mu E(T) + \mu^2 \\
 &= E(T^2) - \mu^2 \\
 &= E(T^2) - [E(T)]^2
 \end{aligned}$$

where T is a dummy variable representing the values which an r_{ik} can take on. The area above the cumulative distribution versus T curve, computed on the basis of horizontal strips, is the integral of the height of the strip, which is T , times the width of the strip, which is $dF(T)$. But

$$dF(T) = f(T)dT$$

and therefore,

$$\text{Area}_T = \int_0^1 T f(T)dT$$

The expected value of T is, by definition,

$$E(T) = \int_0^1 T f(T) dT$$

The area above the cumulative distribution versus T^2 curve is, similarly,

$$\text{Area}_{T^2} = \int_0^1 T^2 f(T) dT$$

and, by definition,

$$E(T^2) = \int_0^1 T^2 f(T) dT$$

Therefore

$$\begin{aligned} \text{Area}_{T^2} - (\text{Area}_T)^2 &= E(T^2) - [E(T)]^2 \\ &= \text{var}(T) \end{aligned}$$

As mentioned previously, the plots of the pseudo-normal, cumulative distributions matched satisfactorily, visually, the plots of the simulated cumulative distributions, except for the late age intervals. It should be noted, perhaps, that the Poisson cumulative distributions were generated for the late age intervals and that they matched satisfactorily, visually, the corresponding simulated cumulative distributions. This result was not utilized in the development of the polynomial fitting procedure.

A computer program for fitting a polynomial to a set of observed values by the weighted or unweighted least-squares method was obtained from the Iowa State University Statistical Laboratory - Numerical Analysis and Programming Section, Ames, Iowa. The program uses orthogonal polynomials to obtain the polynomial coefficients of the least-squares fit. The program requires (1) that the values of the independent variable be equally spaced and (2) that only a single value of the dependent variable be paired with a single value of the independent variable.

The first requirement prohibits the use of the retirement ratio(s) for the age interval 0 to 0.5 years. If there are multiple values of the dependent variable for each value of the independent variable, the second requirement forces the programmer to combine such multiple values into a single value before using the orthogonal polynomial program.

The method of orthogonal polynomials facilitates testing the significance of each additional degree. The program obtained from the Statistical Laboratory computed, and printed out, the regression sum of squares, the remainder sum of squares, the regression mean square, the remainder mean square, the total sum of squares, and the degrees of freedom associated with each. Thus, an F test of the form

$$F_{v_1, v_2, \alpha} = \frac{R_1 - R_2}{R_2/v_2}$$

R_1 = remainder sum of squares of the $(n - 1)^{\text{th}}$ degree polynomial

R_2 = remainder sum of squares of the n^{th} degree polynomial

v_1 = degrees of freedom of $(R_1 - R_2)$
 $= 1$

v_2 = degrees of freedom of R_2
 $= (\text{number of observed values}) - n - 1$

α = probability of a type I error

can readily be performed to test the significance of adding the n^{th} degree term.

An analyst may wish to test the normality of the retirement ratios (at an age interval) obtained from historical property data. A method of

testing the hypothesis

$$\{Y_i\} \sim N(\mu, \sigma) \quad i = 1, 2, \dots, I$$

$\{Y_i\}$ = The ordered set of observed values of a sample (such as the retirement ratios at a particular age interval)

was developed and is presented in Appendix F. A severely limiting condition on the application of the method to retirement ratios is that the sample of retirement ratios to be tested must have come from vintage groups of the same size.

If some r_{ik} is of the form zero divided by zero, the computer program to implement the procedure does not attempt to compute the corresponding $\hat{\sigma}_{ik}^2$ (step 5 of the procedure); the part of the program which calculates $\hat{\sigma}_{ik}^2$ skips that age interval (and all subsequent age intervals of the same vintage group) and proceeds to start calculating the $\hat{\sigma}_{ik}^2$ for the next vintage group. In essence, this process assigns a value of zero to both r_{ik} and w_{ik} when r_{ik} is of the form zero divided by zero.

An interpolated $\tilde{r}_{.k}$ (step 4 of the procedure) may be equal to zero (i.e., a retirement ratio of the form zero divided by a positive constant). The variance of the r_{ik} , $\hat{\sigma}_{ik}^2$, is theoretically zero when $\tilde{r}_{.k}$ is zero and the weights $w_{1k}, w_{2k}, \dots, w_{Ik}$ approach infinity. A digital computer cannot calculate the value of $1/0$; hence, the computer was programmed to print out a code indicating that an $\tilde{r}_{.k}$ of the form zero divided by a constant had been encountered and then to proceed to computations involving $\tilde{r}_{.k} + 1$. Some arbitrary value must be assigned to each of the w_{ik} ($i = 1, \dots, I$). A possible alternative would be to assign a value (to such w_{ik}) which is equal to the largest value of any other w_{ik} , or perhaps a value up to twice as large as the largest value of any other w_{ik} . The

reason for recommending this alternative is to prevent one (or a limited number of) r_{ik} value from dominating or inordinately influencing the calculation of the polynomial coefficients in step (6). No study has been made of what might be the relative magnitude of an appropriate value to assign to such w_{ik} .

DISCUSSION

Two slightly different procedures for fitting polynomials to retirement ratios are currently in use. A third procedure is presented as a preliminary approach in Appendix D and is utilized in the procedure developed in this dissertation. When there is more than one retirement ratio for an age interval, because of the method used to obtain the original life table, all four of these procedures combine the several retirement ratios for an age interval into some single, composite retirement ratio. Three of the procedures then fit a polynomial to the retirement ratios by minimizing the sum of the "weighted" squares of the deviations of the composite retirement ratios from the regression curve; the other procedure gives each squared deviation equal weight.

To facilitate referring to the various procedures, the following designations will be adopted:

1. Procedure A, a currently used method, is

$$\text{Min}_{a,b,c,\text{etc.}} \left\{ \sum_{k=1}^K [\tilde{r}_{\cdot k} - (a + bx_k + cx_k^2 + \dots)]^2 \right\}$$

where

$$\begin{aligned} \tilde{r}_{\cdot k} &= \text{a composite retirement ratio for the } k^{\text{th}} \text{ age interval} \\ &= \frac{\sum_{i=1}^I \tilde{w}_{ik} r_{ik}}{\sum_{i=1}^I \tilde{w}_{ik}} \end{aligned}$$

r_{ik} = retirement ratio for the k^{th} age interval from the i^{th} vintage group

$$= \frac{R_{ik}}{S_{ik}}$$

R_{ik} = number of units from the i^{th} vintage group retired during the k^{th} age interval

S_{ik} = number of units from the i^{th} vintage group surviving at the beginning of the k^{th} age interval

$$\tilde{w}_{ik} = 1/\tilde{\sigma}_{ik}^2$$

$\tilde{\sigma}_{ik}^2$ = the conditional variance of r_{ik} (conditional upon the denominator of r_{ik} , S_{ik})

$$= \frac{P_k Q_k}{S_{ik}}$$

P_k = population retirement ratio for the k^{th} age interval

$$Q_k = 1 - P_k$$

Therefore

$$\begin{aligned} \tilde{r}_{\cdot k} &= \frac{\sum_{i=1}^I (S_{ik}/P_k Q_k) r_{ik}}{\sum_{i=1}^I S_{ik}/P_k Q_k} \\ &= \frac{\sum_{i=1}^I S_{ik} r_{ik}}{\sum_{i=1}^I S_{ik}} \end{aligned}$$

since P_k and Q_k are constants for a given k and $i = 1, 2, \dots, I$.

Also

$a, b, c, \text{etc.}$ = polynomial coefficients

$x_k = k^{\text{th}}$ age interval index number

2. Procedure B, a currently used method, is

$$\text{Min}_{a, b, c, \text{etc.}} \left\{ \sum_{k=1}^K S_{\cdot k} [\tilde{r}_{\cdot k} - (a + bx_k + cx_k^2 + \dots)]^2 \right\}$$

where

$$S_{\cdot k} = \sum^I S_{ik}$$

$\tilde{r}_{\cdot k}$, a , b , c , and x_k are as shown in (1), above.

3. Procedure C, the preliminary approach suggested in Appendix D, is

$$\text{Min}_{a,b,c,\text{etc.}} \left\{ \sum^K \tilde{w}_{\cdot k} [\tilde{r}_{\cdot k} - (a + bx_k + cx_k^2 + \dots)]^2 \right\}$$

where

$$\begin{aligned} \tilde{w}_{\cdot k} &= \frac{S_{\cdot k}}{\tilde{r}_{\cdot k}(1 - \tilde{r}_{\cdot k})} \\ &= 1/\tilde{\sigma}_{\cdot k}^2 \end{aligned}$$

$$\begin{aligned} \tilde{\sigma}_{\cdot k}^2 &= \text{var}(\tilde{r}_{\cdot k} | S_{\cdot k}) \\ &= \frac{P_k Q_k}{S_{\cdot k}} \end{aligned}$$

and perhaps, is best described as a weighted sum of conditional variances. P_k and Q_k must be replaced by their sample estimates, $\tilde{r}_{\cdot k}$ and $(1 - \tilde{r}_{\cdot k})$, respectively, since they are not known. Hence

$$\tilde{\sigma}_{\cdot k}^2 = \frac{\tilde{r}_{\cdot k}(1 - \tilde{r}_{\cdot k})}{S_{\cdot k}}$$

$S_{\cdot k}$, $\tilde{r}_{\cdot k}$, a , b , c , and x_k are as shown in (1) and (2), above.

4. Procedure D, the procedure developed in the preceding section, is

$$\text{Min}_{a,b,c,\text{etc.}} \left\{ \sum^K w_{\cdot k} [r_{\cdot k} - (a + bx_k + cx_k^2 + \dots)]^2 \right\}$$

where

$$w_{\cdot k} = \sum^I 1/\sigma_{ik}^2$$

σ_{ik}^2 = estimate of the variance of r_{ik} obtained by means of the

pseudo-normal, cumulative distribution computer program and based on (1) C_k and C'_k values calculated from the preliminary fit of a polynomial to the retirement ratios and (2) J_i

$$r_{\cdot k} = \frac{\sum_{i=1}^I w_{ik} r_{ik}}{\sum_{i=1}^I w_{ik}}$$

$$w_{ik} = 1/\sigma_{ik}^2$$

a , b , c , and x_k are as shown in (1), above.

Three basic assumptions of all four procedures are:

1. The r_{ik} are independent random samples,
2. The age intervals during which the units are retired are determined without error, and
3. The $E(r_{ik})$ is constant for a given k and $i = 1, 2, \dots, I$.

From a practical point of view:

1. The first assumption appears to be valid for a given k and $i = 1, 2, \dots, I$ but not valid for a given i and $k = 1, 2, \dots, K$.
2. The second assumption is probably valid.
3. The third assumption appears to be invalid, in general.

The four procedures will be discussed in this section in terms of the properties (unbiasedness and minimum variance) of:

1. The estimators of the composite retirement ratios, and
2. The estimators of the polynomial coefficients.

An additional topic which will be briefly discussed in this section is the effect of using dollars rather than physical units (as the measure of the

amount of property) on calculating estimates of the variances.

The criteria for evaluating each procedure, in a qualitative sense, are:

1. Does the procedure utilize unbiased estimators and, if so, in what sense are these estimators unbiased?
2. Do these estimators have any good variance properties?
3. Do the estimators coincide with the estimators obtained by the principle of least-squares?

Estimators of the Composite Retirement Ratios

Procedures A, B, and C combine the r_{ik} for an age interval into a single composite retirement ratio, $\tilde{r}_{.k}$, in the same manner.

$$\begin{aligned}\tilde{r}_{.k} &= \frac{\sum_{i=1}^I (S_{ik}/P_k Q_k) r_{ik}}{\sum_{i=1}^I S_{ik}/P_k Q_k} \\ &= \frac{\sum_{i=1}^I S_{ik} r_{ik}}{\sum_{i=1}^I S_{ik}}\end{aligned}$$

Since:

1. Each r_{ik} is assumed to be an unbiased estimator of P_k , and
 2. The conditional variance of each r_{ik} can be assumed to be known because S_{ik} is known and P_k and Q_k cancel out,
- the $\tilde{r}_{.k}$ could be said to be linear unbiased estimators of P_k having the minimum variance of all linear unbiased estimators. Graybill provides the necessary theorem (9, p. 409):

Theorem 18.11. Let $\hat{\theta}_1$ be an unbiased estimator of θ , and let the variance of $\hat{\theta}_1$ be denoted by σ_1^2 . Let $\hat{\theta}_2$ be another unbiased estimator of θ , and let the variance of $\hat{\theta}_2$ be denoted by σ_2^2 . Let $\hat{\theta}_1$ and $\hat{\theta}_2$ be uncorrelated. Then the best (minimum-variance) linear unbiased estimator of θ is

$$\hat{\theta} = \frac{\sigma_2^2 \hat{\theta}_1 + \sigma_1^2 \hat{\theta}_2}{\sigma_1^2 + \sigma_2^2}$$

Repeated application of Theorem 18.11 yields

$$\hat{\theta} = \frac{\sum_{i=1}^I (1/\sigma_i^2) \hat{\theta}_i}{\sum_{i=1}^I 1/\sigma_i^2}$$

Hence, on the assumption that the conditional variances are the correct variances, the $\tilde{P}_{\cdot k}$ are best linear unbiased estimators of the P_k . When S_{ik} is large and r_{ij} small, the conditional variances will very closely approximate the correct variances. The least-squares estimator of P_k is (11, p. 12)

$$\hat{\theta} = \frac{\sum_{i=1}^I w_{ik} r_{ik}}{\sum_{i=1}^I w_{ik}}$$

where

$$w_{ik} = \sigma^2 / \sigma_{ik}^2$$

σ^2 = a constant of proportionality which may be used to alter the magnitudes of the w_{ik}

Hence, substituting $\tilde{\sigma}_{ik}^2$ for σ_{ik}^2

$$\hat{\theta} = \frac{\sum_{i=1}^I (\sigma^2 / \tilde{\sigma}_{ik}^2) r_{ik}}{\sum_{i=1}^I \sigma^2 / \tilde{\sigma}_{ik}^2}$$

$$= \frac{\sum \frac{1}{\hat{\sigma}_{ik}^2} r_{ik}}{\sum \frac{1}{\hat{\sigma}_{ik}^2}}$$

and, therefore, $\tilde{r}_{\cdot k}$ is of the same form as the least-squares estimator of P_k with the conditional variances used as the variances of the r_{ik} .

Procedure D utilizes a preliminary polynomial fit of the retirement ratios (by procedure C) to obtain "better" estimates of the C_k and C'_k used to calculate $\hat{\sigma}_{ik}^2$. If, because of these refined estimates of C_k and C'_k , one is willing to assume that the $\hat{\sigma}_{ik}^2$ are the actual variances of the r_{ik} , then the $r_{\cdot k}$ are best linear unbiased estimators of the P_k (9, p. 409). When the $\hat{\sigma}_{ik}^2$ are considered as sample estimates of the variances of the r_{ik} , very little can be said about the unbiasedness and variance properties of the $r_{\cdot k}$ as estimators of the P_k . The $r_{\cdot k}$ are of the same form as the least-squares estimators of the P_k regardless of whether the $\hat{\sigma}_{ik}^2$ are considered as the actual variances or as sample estimates of the variances of the r_{ik} . The maximum-likelihood estimators of the P_k , assuming the r_{ik} are distributed $N(\mu_k, \sigma_{ik})$, are identical to the least-squares estimators. If the additional assumption is made that the r_{ik} are distributed $N(\mu_k, \sigma_{ik})$, then each $r_{\cdot k}$ is distributed $N(\mu_k, \sigma_{\cdot k})$ (11, pp. 29-30) where

$$\sigma_{\cdot k}^2 = \frac{1}{\sum \frac{1}{\sigma_{ik}^2}}$$

Estimators of the Polynomial Coefficients

The polynomial fitting portion of procedure A assigns equal weight to the $\tilde{r}_{\cdot k}$. According to Guest (11, pp. 88-89), the estimators of the

polynomial coefficients obtained by solving the set of equations

$$\sum^K \lambda_k [Y_k - (a + bx_k + cx_k^2 + \dots)] x_k = 0$$

where

$$\{\lambda_k\} = \text{set of known constants}$$

will be unbiased. Therefore procedure A does yield unbiased estimators of the polynomial coefficients. When the variances, $\sigma_{\cdot k}^2$, are not equal to a constant, then the weights should be

$$\{\lambda_k\} = \{1/\sigma_{\cdot k}^2\}$$

Hence, the estimators of the polynomial coefficient do not have good variance properties. The estimators of the polynomial coefficients are not the same as the least-squares estimators since procedure A sets λ_k equal to a constant and the principle of least-squares sets λ_k equal to $1/\sigma_{\cdot k}^2$.

Procedure B weights each $\tilde{r}_{\cdot k}$ by the corresponding $S_{\cdot k}$. The inverse of $S_{\cdot k}$, $1/S_{\cdot k}$, might be considered as an estimate of the variance of $\tilde{r}_{\cdot k}$, but it is not the best available estimate; $1/S_{\cdot k}$ is not even a correct estimate of the conditional variance. The appropriate estimator of the variance $\tilde{\sigma}_{\cdot k}^2$, consistent with the variance $\tilde{\sigma}_{ik}^2$, is

$$\begin{aligned} \tilde{\sigma}_{\cdot k}^2 &= \text{var}(\tilde{r}_{\cdot k} | S_{ik}) \\ &= \frac{P_k Q_k}{S_{\cdot k}} \end{aligned}$$

Procedure B utilizes unbiased estimators of the polynomial coefficients if the S_{ik} are assumed to be known constants in accordance with the conditional variance assumption. These estimators should have somewhat better variance properties than the estimators utilized in procedure A because $1/S_{\cdot k}$ is a

better approximation of $\text{var}(\tilde{r}_{\cdot k} | S_{ik})$ than the constant used in procedure A. The estimators of the polynomial coefficients are not the same as the least-squares estimators to the extent that $1/S_{\cdot k}$ is not the same as $\text{var}(\tilde{r}_{\cdot k} | S_{ik})$.

The weights used in procedure C to obtain the estimators of the polynomial coefficients are at least in agreement with the conditional variance assumption.

$$\begin{aligned}\tilde{w}_{\cdot k} &= 1/\text{var}(r_{\cdot k} | S_{ik}) \\ &= \frac{S_{\cdot k}}{P_k Q_k} \\ &= \frac{S_{\cdot k}}{\tilde{r}_{\cdot k} (1 - \tilde{r}_{\cdot k})}\end{aligned}$$

$\tilde{w}_{\cdot k}$ is the inverse of the sample estimate of the variance of $\tilde{r}_{\cdot k}$ because P_k and Q_k are not known and must be replaced by their sample estimates $\tilde{r}_{\cdot k}$ and $(1 - \tilde{r}_{\cdot k})$. The estimators of the polynomial coefficients are not necessarily unbiased since the $\tilde{w}_{\cdot k}$, where

$$\{\tilde{w}_{\cdot k}\} = \{\lambda_k\}$$

are not known constants. These estimators should have, perhaps, somewhat better variance properties than the estimators of either procedure A or B. These estimators (procedure C) are of the same form as the least-squares estimators since

$$\begin{aligned}\{\lambda_k\} &= \{\tilde{w}_{\cdot k}\} \\ &= \{1/\tilde{\sigma}_{\cdot k}^2\}\end{aligned}$$

Procedure D yields best linear unbiased estimators of the polynomial coefficients if the $\tilde{\sigma}_{ik}^2$ are assumed to be the actual variances of the r_{ik} .

$$\begin{aligned}
\text{var}(r_{\cdot k}) &= \text{var}\left[\frac{\sum_{i=1}^I (1/\hat{\sigma}_{ik}^2)(r_{ik})}{\sum_{i=1}^I 1/\hat{\sigma}_{ik}^2}\right] \\
&= \left(\frac{1}{\sum_{i=1}^I 1/\hat{\sigma}_{ik}^2}\right)^2 [\text{var}(r_{1k}/\hat{\sigma}_{1k}^2 + r_{2k}/\sigma_{2k}^2 + \dots + r_{Ik}/\sigma_{Ik}^2)] \\
&= \left(\frac{1}{\sum_{i=1}^I 1/\hat{\sigma}_{ik}^2}\right)^2 \left[\left(\frac{1}{\hat{\sigma}_{ik}^2}\right)^2 \text{var}(r_{1k}) + \left(\frac{1}{\hat{\sigma}_{2k}^2}\right)^2 \text{var}(r_{2k}) + \dots \right. \\
&\quad \left. + \left(\frac{1}{\hat{\sigma}_{Ik}^2}\right)^2 \text{var}(r_{Ik}) \right] \\
&= \left(\frac{1}{\sum_{i=1}^I 1/\hat{\sigma}_{ik}^2}\right)^2 \left[\left(\frac{1}{\hat{\sigma}_{ik}^2}\right)^2 \hat{\sigma}_{1k}^2 + \left(\frac{1}{\hat{\sigma}_{2k}^2}\right)^2 \sigma_{2k}^2 + \dots \right. \\
&\quad \left. + \left(\frac{1}{\hat{\sigma}_{Ik}^2}\right)^2 \sigma_{Ik}^2 \right] \\
&= \left(\frac{1}{\sum_{i=1}^I 1/\hat{\sigma}_{ik}^2}\right)^2 \sum_{i=1}^I 1/\sigma_{ik}^2 \\
&= \frac{1}{\sum_{i=1}^I 1/\hat{\sigma}_{ik}^2} \\
&= 1/w_{\cdot k}
\end{aligned}$$

Under the (additional) assumption of normality, the estimators are unbiased and have minimum variance amongst unbiased estimators (9, p. 117; 11, pp. 88-89). If the assumption

$$\hat{\sigma}_{ik}^2 = \text{var}(r_{ik})$$

is not made then the estimators of the polynomial coefficients are not necessarily unbiased and do not have optimum variance properties. The form of these estimators is the same as the form of the least-squares estimators.

In summary:

1. Procedures A, B, and C utilize best linear unbiased estimators of the composite retirement ratios if the conditional variances $\hat{\sigma}_{ik}^2$ are assumed to be the actual variances of the r_{ik} .
2. Procedure D utilizes best linear unbiased estimators of the composite retirement ratios if the $\hat{\sigma}_{ik}^2$ are assumed to be the actual variances of the r_{ik} .
3. Procedure A utilizes unbiased estimators of the polynomial coefficients; these estimators are not of the same form as the appropriate least-squares estimators and probably have relatively poor variance properties.
4. Procedure B utilizes unbiased estimators of the polynomial coefficients if the conditional variance assumption is made; these estimators are not of the same form as the least-squares estimators, but should have somewhat better variance properties than the procedure A estimators.
5. The estimators of the polynomial coefficients utilized in procedure C are not necessarily unbiased but should have somewhat better variance properties than the procedure A and procedure B estimators; these estimators are of the same form as the least-squares estimators.

6. Procedure D utilizes best linear unbiased estimators of the polynomial coefficients if the $\hat{\sigma}_{ik}^2$ are assumed to be the actual variances of the r_{ik} ; these estimators are of the same form as the least-squares estimators.
7. If the r_{ik} are assumed to be distributed $N(\mu_k, \sigma_{ik})$ and if $\hat{\sigma}_{ik}^2$ is assumed to be the actual variance, σ_{ik}^2 , then the estimators of the polynomial coefficients utilized in procedure D are unbiased and have minimum variance amongst all unbiased estimators.

Procedure C utilizes weights of $\tilde{w}_{\cdot k}$ based on the variances $\hat{\sigma}_{\cdot k}^2$ in fitting polynomials to the $\tilde{r}_{\cdot k}$, as previously mentioned. These variances are conditional upon the denominators, S_{ik} , of the retirement ratios and, when all S_{ik} are known, are variances of constants. Theoretically the variance of a constant is zero. The extent to which this theoretical consideration limits the usefulness of procedure C is not known.

The difference between the results that might be obtained from the practical application of procedures C and D is not known. Either C or D should yield better results than either A or B.

Effect of Dollars on Computing Variances

The primary effect of using dollars rather than physical units on variances, assuming that each physical unit costs more than one dollar, is to decrease the magnitude of the variances. The problem which this effect engenders is the calculation of extremely small variances. This problem does not seem to arise in procedures A, B, and C since:

1. The variances $\hat{\sigma}_{ik}^2$ are never directly computed, and

2. The variances $S_{.k}$ and $\hat{\sigma}_{.k}^2$ are calculated directly from the raw data.

The variance estimate used in Appendix D, $\hat{\sigma}_{ik}^2$, is based upon the calculation of Area_T (the area above the pseudo-normal, cumulative distribution versus T curve) and Area_{T^2} (the area above the pseudo-normal, cumulative distribution versus T^2 curve). These areas, in turn, are based upon the cumulative distribution generated by the pseudo-normal computer program which utilizes the parameters: vintage group size, C_k and C'_k . As the vintage group size increases, the variance $\hat{\sigma}_{ik}^2$ decreases, and therefore the amount by which T is incremented must also decrease or the variance will appear to be zero. For instance, a weight (of the form $S_{.k}$) of 1,000,000 corresponds to a variance of 0.000,001. Hence, to obtain a $\hat{\sigma}_{ik}^2$ of the same magnitude, T should be incremented by amounts of approximately 0.000,000,1 in order to obtain any accuracy in the estimate of $\hat{\sigma}_{ik}^2$. The problem then becomes one of computer time and the quick determination of the proper amount by which to increment T and of the first non-zero value of

$$\Pr\{N[(1-T)J C_k - T J C'_k, \{(1-T)^2 J C_k(1-C_k) + T^2 J C'_k(1-C'_k) + 2T(1-T)J C_k C'_k\}^{1/2}] \leq 0\}$$

This problem does not appear to be insolvable.

EXAMPLE

The procedure developed in Section VI was applied to a four vintage group example. The amount of property in each vintage group was measured in physical units and was determined by drawing a random uniform number between 150 and 500; the sizes of the vintage groups were 462, 176, 348, and 226 units. The retirement experience of each vintage group was simulated on the basis of a R_1 type Iowa curve and an average service life of 25 years (a flow chart of the simulation computer program is presented in Appendix A).

A plot of the four retirement ratios at each age interval (three at the last age interval) is shown in Figures 13a and 13b. A larger ordinate scale is used in Figure 13a than in Figure 13b to avoid crowding the points close together. Age interval one is the age interval 0 to 0.5 years, age interval two is the age interval 0.5 to 1.5 years, etc.

One of the retirement ratio methods of analyzing historical data is to fit polynomials of up to the fourth degree to retirement ratios of the form

$$\hat{r}_{\cdot k} = \frac{L_{\cdot k}}{L_{\cdot k} + M_{\cdot k}}$$

by the unweighted least-squares method. The polynomial selected as best representing the original data is generally the first, second, or third degree polynomial. Plots of the second and third degree polynomials and of the ninth degree polynomial (which is the highest degree polynomial significant at the 0.05 level by the F test) are shown in Figures 14 and 15, respectively.

Figure 13a. Vintage group retirement ratios at each age interval for the example

- - vintage group I (462 units)
- - vintage group II (176 units)
- △ - vintage group III (348 units)
- ◇ - vintage group IV (226 units)

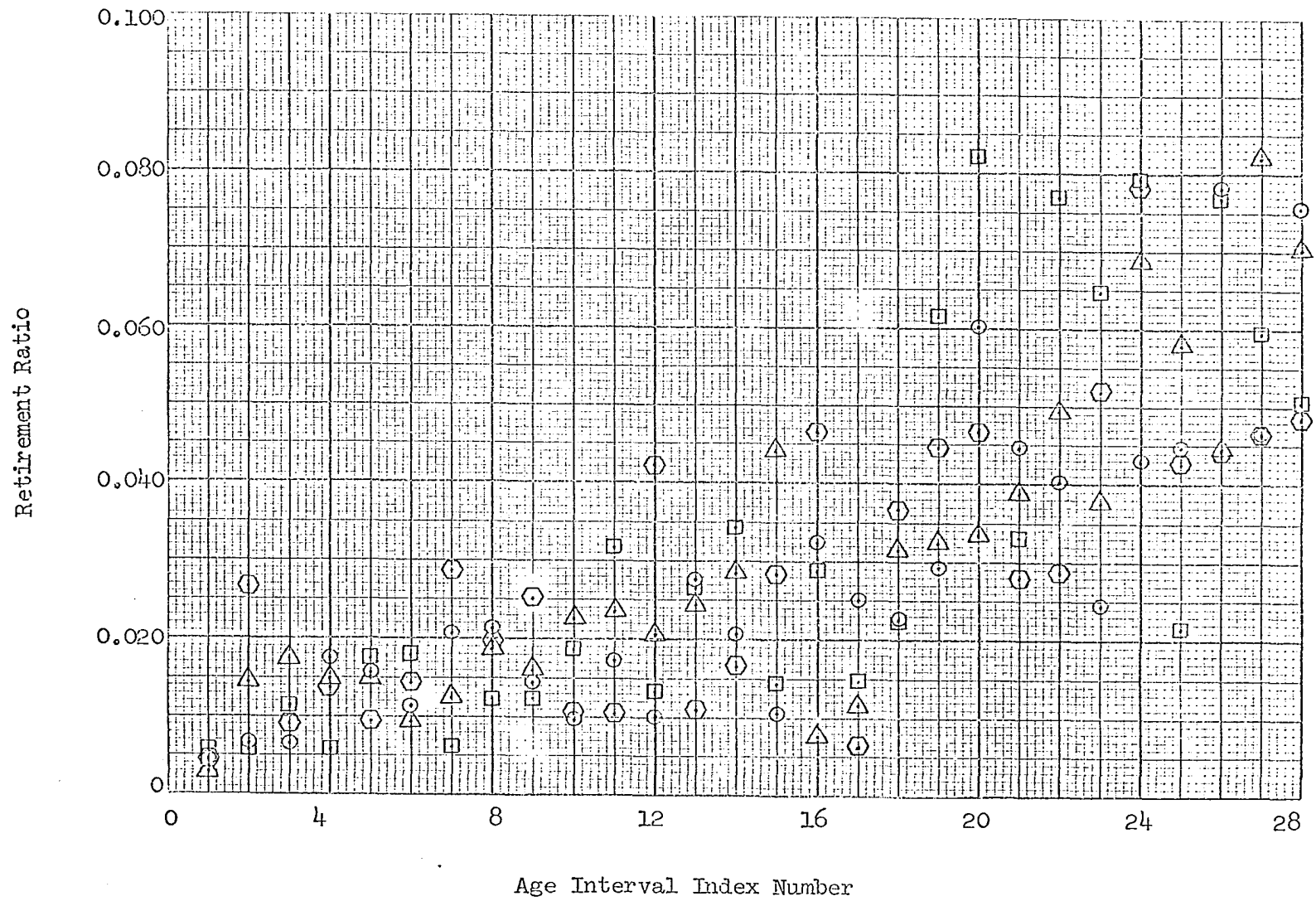


Figure 13b. Vintage group retirement ratios at each age interval for the example

- - vintage group I (462 units)
- - vintage group II (176 units)
- △ - vintage group III (348 units)
- - vintage group IV (226 units)

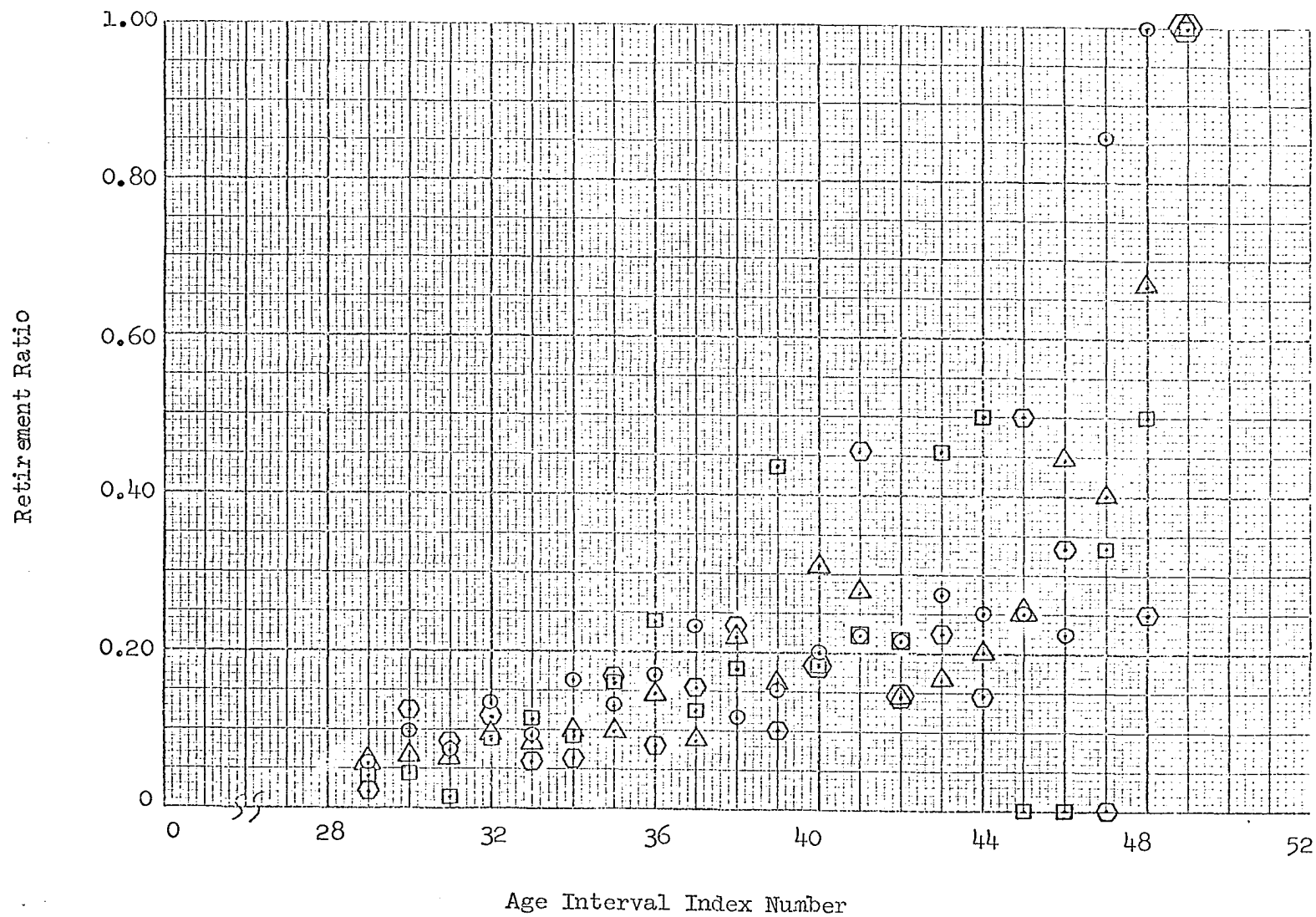


Figure 14. Second and third degree polynomial fits of the $\tilde{r}_{\cdot k}$

○ - second degree polynomial

□ - third degree polynomial

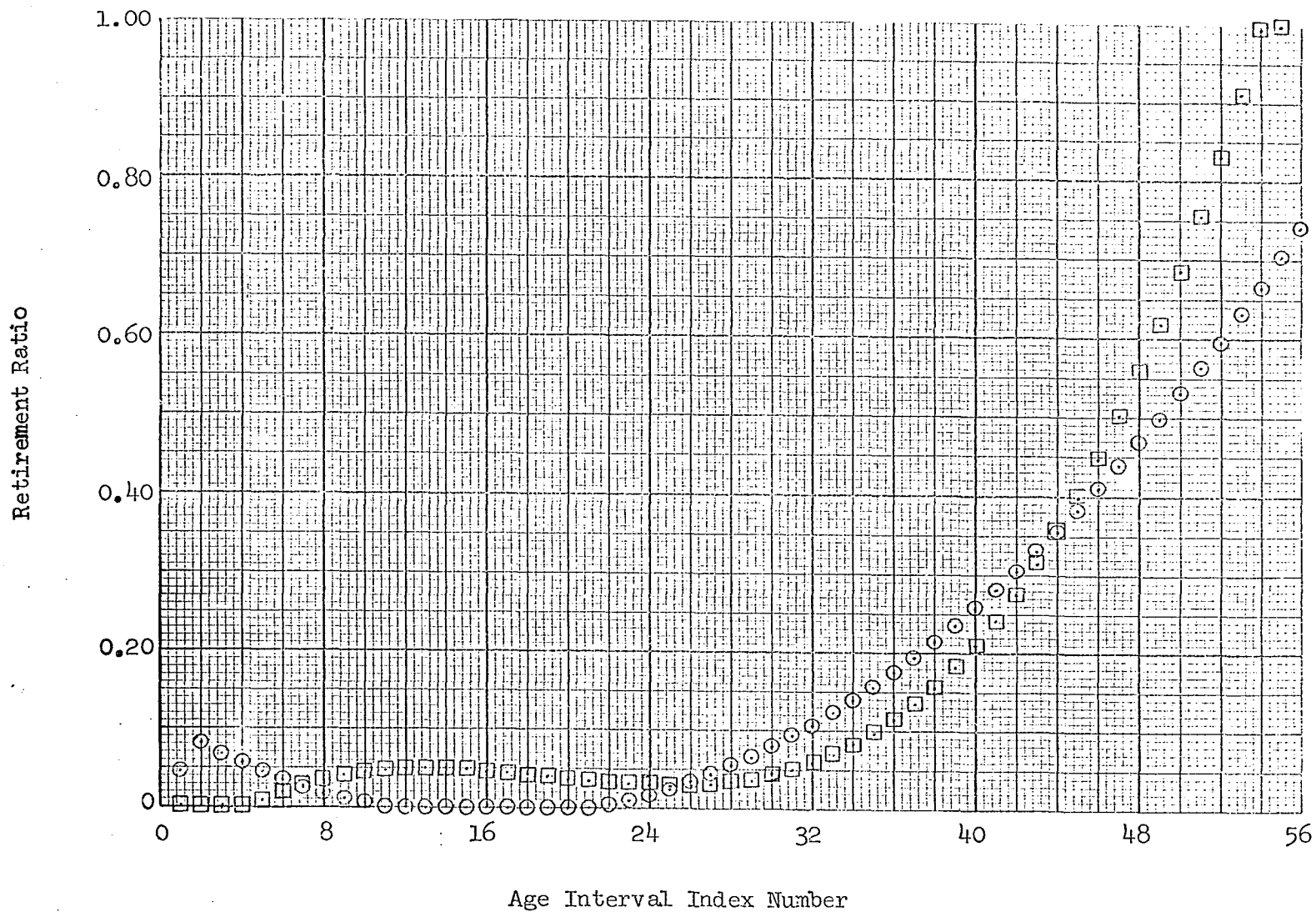
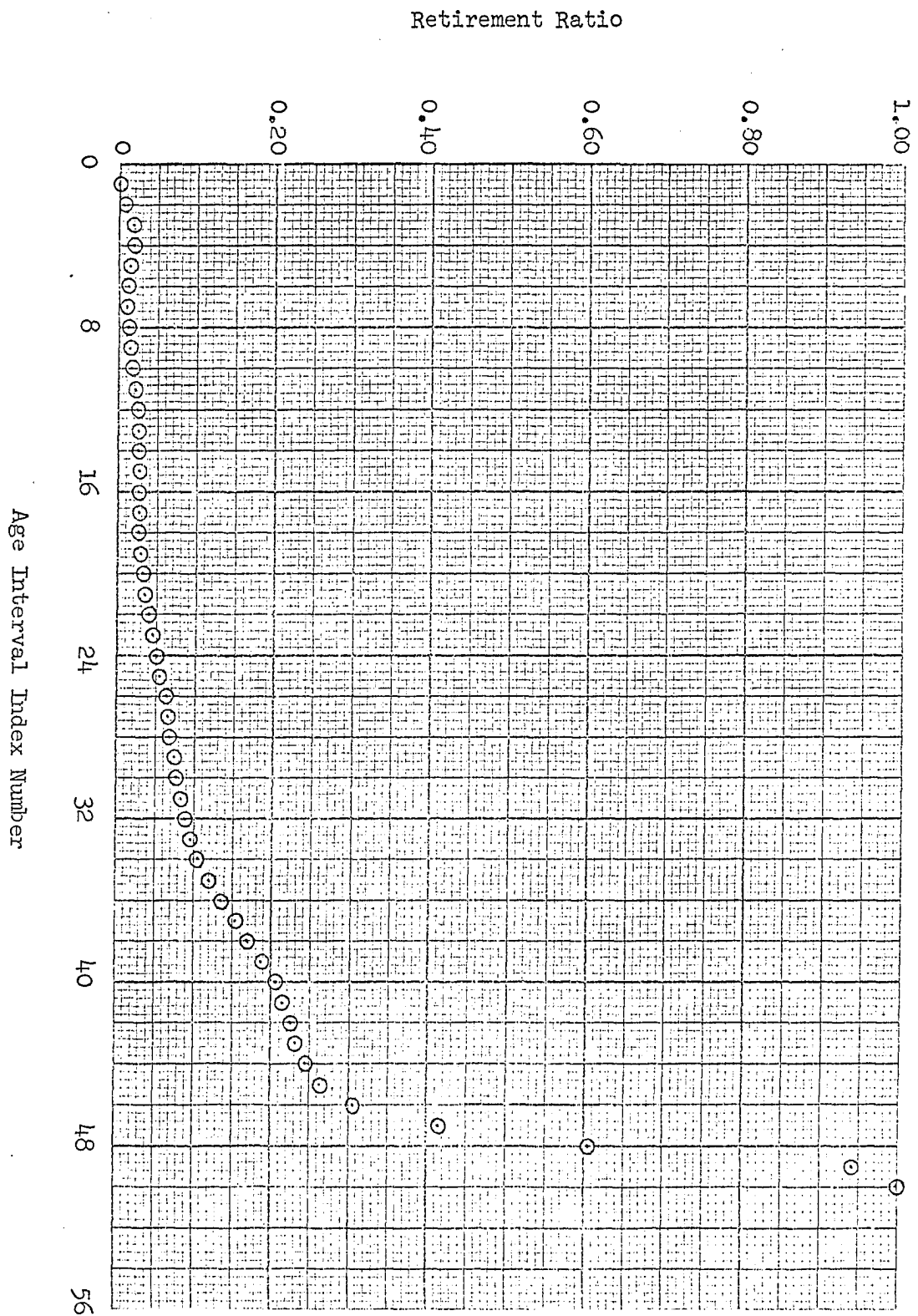


Figure 15. Ninth degree polynomial fit of the $\tilde{r}_{\bullet k}$



The orthogonal polynomial program, obtained from the Statistical Laboratory, was used to fit the polynomials to the retirement ratios. This program requires equal spacing of the abscissa values, hence, the retirement ratio for the age interval 0 to 0.5 years was not used as an input value. All other age intervals were assigned an index number one less than was previously assigned to them, so that age interval one represented the age interval 0.5 to 1.5 years, etc. After this reassignment of index numbers, each index number (k) was, numerically, midway between the boundary years of age of the age interval it represented (i.e., a k value of one represented age interval 0.5 to 1.5 years, etc.). The value of the retirement ratio for the age interval 0 to 0.5 years was extrapolated from the polynomial as one-half of the retirement ratio for a k value of 0.25, a k value which is, numerically, midway between 0 years and 0.5 years.

The age intervals were then reassigned their previous index numbers for plotting purposes, so that in the figures, an index number of one represents the age interval 0 to 0.5 years, an index number of two represents the age interval 0.5 to 1.5 years, etc.

The second degree polynomial does not fit the $\tilde{r}_{.k}$ very satisfactorily. The retirement ratio values interpolated from the polynomial are too large during the early age intervals and are zero from age interval eleven to age interval twenty-one (these "zero" valued retirement ratios were actually negative but were set equal to zero since, for practical purposes, the amount of property cannot increase as the age interval index number increases). The retirement ratios appear to be too small during the late age intervals; however, the calculated percent surviving at age 50.5 years (the age at the end of the maximum age interval of the theoretical

R_1 - 25 curve) was only 0.019%. The survivor curve calculated from the second degree polynomial (not shown) drops sharply from 100% at age 0 to 66.2% at age 9.5 years, is a horizontal line from age 9.5 years to age 20.5 years, and then drops sharply again.

The third degree polynomial fits the $\tilde{F}_{.k}$ somewhat better than the second degree polynomial. However, the interpolated retirement ratios are zero for the first few age intervals (they were actually negative but set equal to zero) and are a little too high for the age intervals 32 to 47. The calculated percent surviving at age 50.5 years was 0.006%. The corresponding, smoothed survivor curve is a horizontal line (at 100%) from age 0 to age 3.5 years and then drops fairly sharply from age 3.5 years to approximately age 16.5 years, drops relatively less sharply from age 16.5 years to age 28.5 years, drops fairly sharply from age 25.5 years to age 41.5 years, and drops less sharply to approximately zero percent surviving at age 49.5 years.

The ninth degree polynomial fits the $\tilde{F}_{.k}$ satisfactorily. The corresponding smooth survivor curve (not shown) is somewhat irregular from age 0 years to age 6.5 years but essentially follows an R_1 type curve, with an average service life of approximately 24.5 years, beyond age 6.5 years. The calculated percent surviving at age 48.5 was 0.003%.

The procedure developed in this dissertation was applied to the data of the example. Polynomials of the first through tenth degrees were fitted to the $\tilde{F}_{.k}$ by the weighted least squares method, where

$$\tilde{F}_{.k} = \frac{L_{..k}}{L_{..k} + M_{..k}}$$

and the weight given each $\tilde{r}_{\cdot k}$ is

$$w_{\cdot k} = \frac{L_{\cdot\cdot k} + M_{\cdot\cdot k}}{\tilde{r}_{\cdot k}(1 - \tilde{r}_{\cdot k})}$$

The highest degree polynomial which was significant at the 0.05 level by the F test was the third degree polynomial. The $\tilde{r}_{\cdot k}$ were interpolated from the above-mentioned third degree polynomial fit of the $\tilde{\tilde{r}}_{\cdot k}$ and used to calculate the $\hat{\sigma}_{ik}^2$. Then polynomials of degree one through ten were fitted to the r_{ik} by the weighted least-squares method, where the weight assigned to each r_{ik} was $1/\hat{\sigma}_{ik}^2$. The fourth degree polynomial was the highest degree polynomial which was significant at the 0.05 level by the F test.

A plot of the retirement ratios interpolated from the fourth degree polynomial fit of the r_{ik} is shown in Figure 16. The retirement ratio at age interval 59 (57.5 to 58.5 years) is one. The corresponding smoothed survivor curve is shown in Figure 17 and a plot of the $R_1 - 25$ survivor curve is shown in Figure 18. The smoothed survivor curve, plotted according to the commonly used ordinate and abscissa scales, forms a very smooth curve and almost exactly matches an $R_1 - 24.5$. The percent surviving at age 50.5 years was 0.048%.

Figure 16. Smoothed retirement ratio curve from the fourth degree polynomial fit of the r_{ik}

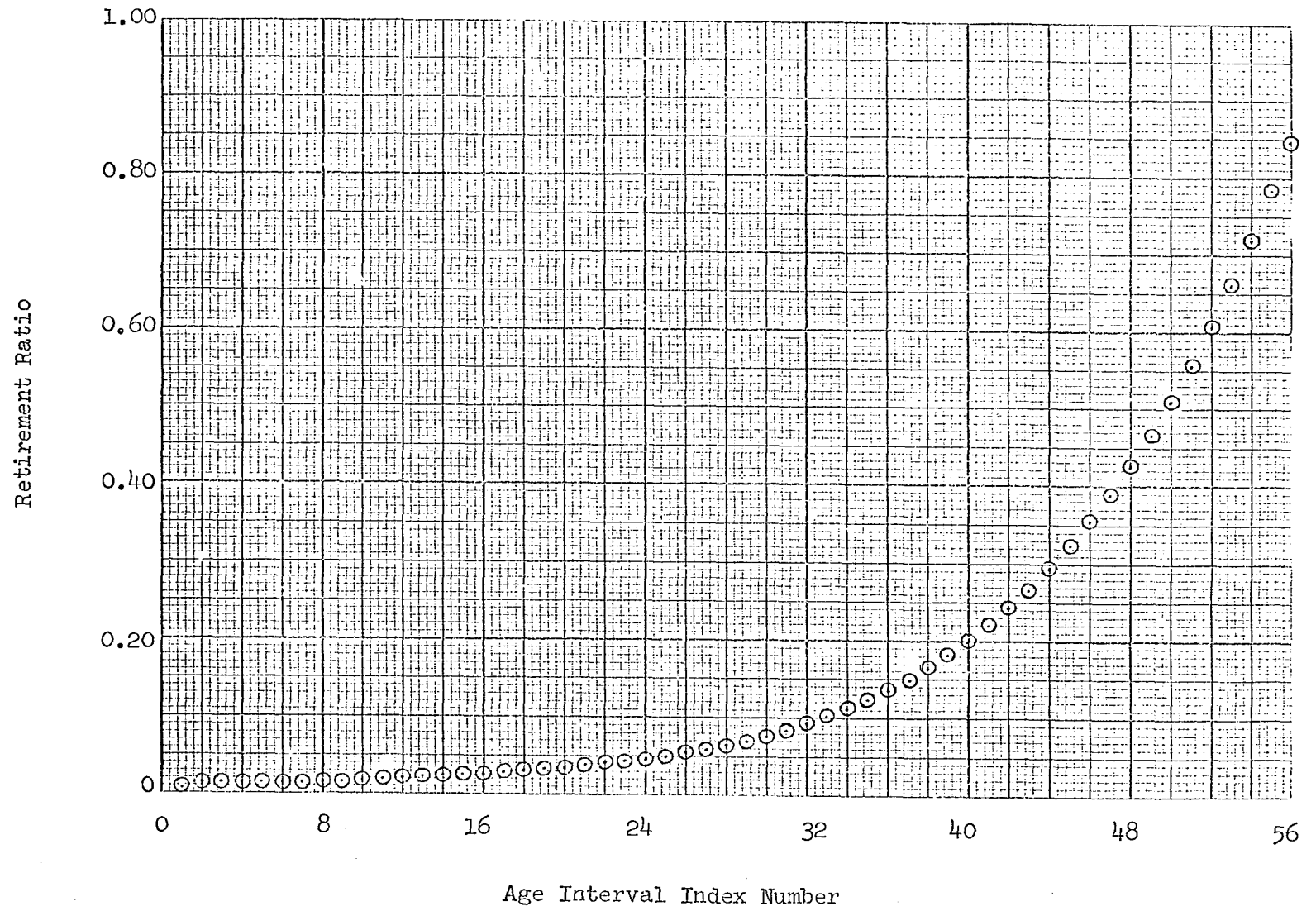


Figure 17. Smoothed survivor curve from the fourth degree, polynomial fit of the r_{ik}

Percent Surviving

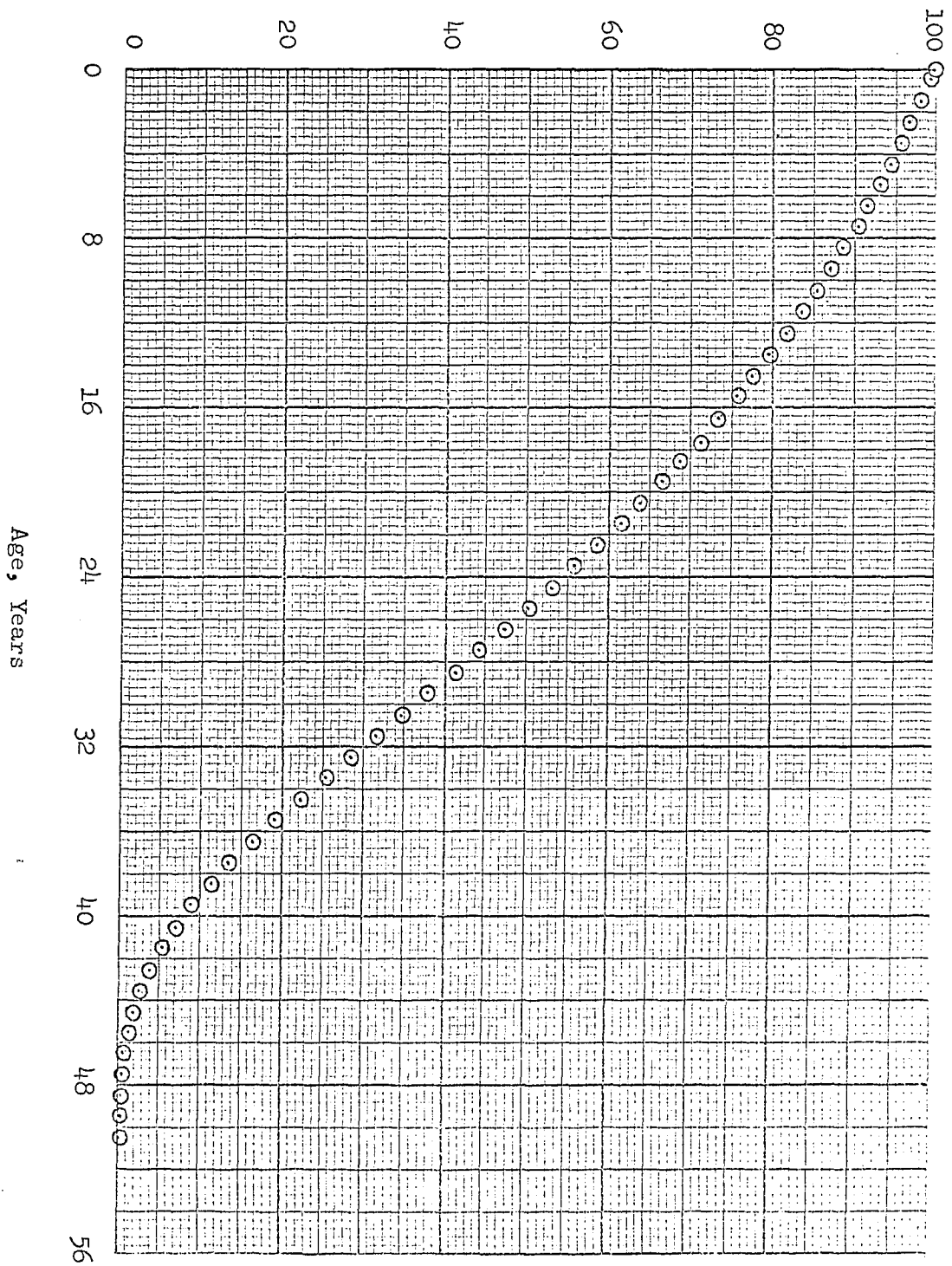
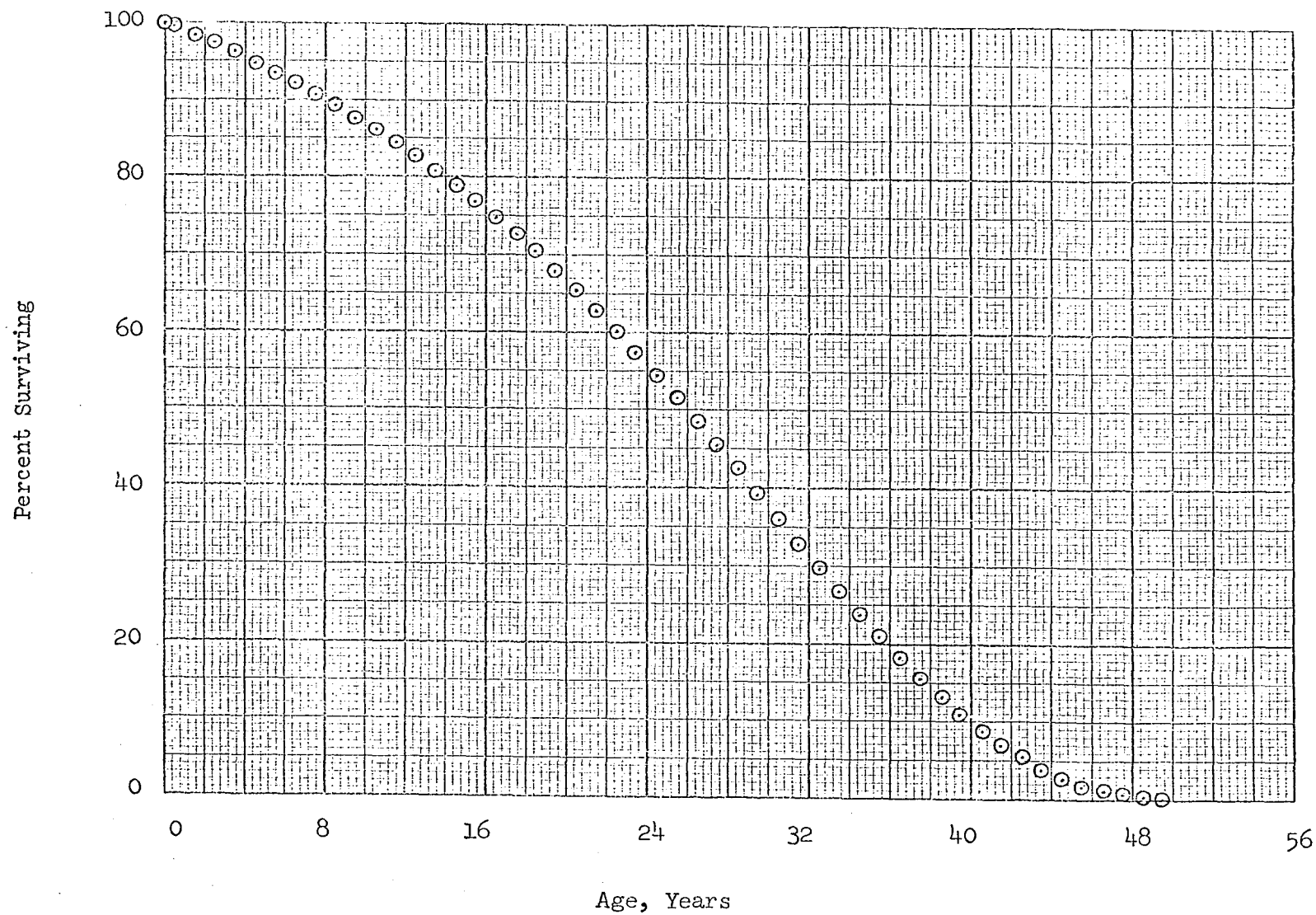


Figure 18. R_1 - 25 survivor curve



CONCLUSIONS

The investigation of the vertical distribution of retirement ratios at each age interval, by simulation, indicated the following:

1. The points of the cumulative distribution of retirement ratios for an age interval plotted on normal probability paper lie nearly along a straight line, except for the early and late age intervals.
2. The variance of the distribution of retirement ratios generally increases as the age interval index number increases (for a given Iowa type curve, average service life, and sample size).
3. The variance of the distribution of retirement ratios for a given age interval decreases as the vintage group size increases (for a given Iowa type curve and average service life).

Hence, each retirement ratio from a vintage group is a sample from an approximately normal distribution and the assumption of the homoscedasticity of variances is invalid.

Further investigation yielded a pseudo-normal computer program which generated cumulative distributions that closely matched, visually, the simulated cumulative distributions of retirement ratios for all age intervals, except the late age intervals. The variance of the cumulative distribution generated by the pseudo-normal computer program can be calculated and only the vintage group size, an estimate of the probability of a unit being retired during the given age interval, and an estimate of the probability of a unit being retired after the given age interval need to be known.

A procedure was developed for fitting polynomials to retirement ratios. The basic assumptions of the procedure are:

1. The r_{ik} are independent random samples,
2. The age intervals during which the units are retired are determined without error, and
3. The $E(r_{ik})$ is a constant for a given k .

Under these assumptions, the procedure utilizes estimators of the polynomial coefficients which are not necessarily unbiased but which probably have relatively good variance properties. If the estimates of the variances of the retirement ratios are assumed to be the actual variances, an assumption which may be reasonable because of the manner in which the variances are calculated, the estimators of the polynomial procedure are best (minimum variance) linear unbiased estimators. If, in addition, the r_{ik} are assumed to be distributed $N(\mu_k, \sigma_{ik})$, the estimators are unbiased and have minimum variance amongst all unbiased estimators; the normality assumption is supported by the approximate linearity of the plots of the simulated cumulative distributions on normal probability paper.

The procedure developed herein was not applied to the data of actual industrial property. Future research needs to be done to determine whether this procedure is significantly better than previously developed procedures and the procedure set forth in Appendix D. Additional research might also be directed towards developing a procedure for fitting polynomials to retirement ratios which considers the effect of the non-independence of the retirement ratios calculated from the same vintage group.

LITERATURE CITED

1. Anderson R. L. and Bancroft, T. A. Statistical theory in research. New York, New York, McGraw-Hill Book Co., Inc. 1952.
2. Bonbright, James C. The valuation of property. Vol. 1. New York, New York, McGraw-Hill Book Co., Inc. 1937.
3. Couch, Frank Van Buskirk, Jr. Classification of type O retirement characteristics of industrial property. Unpublished M.S. thesis. Ames, Iowa, Library, Iowa State University of Science and Technology. 1957.
4. Cowles, Harold Andrews, Jr. Prediction of mortality characteristics of industrial property groups. Unpublished Ph.D. thesis. Ames, Iowa, Library, Iowa State University of Science and Technology. 1957.
5. Edison Electric Institute. Methods of estimating utility plant life. New York, New York, Edison Electric Institute. 1952.
6. Fitch, W. Chester. Fundamental aspects of depreciation theory. Unpublished Ph.D. thesis. Ames, Iowa, Library, Iowa State University of Science and Technology. 1950.
7. Freund, John E. Modern elementary statistics. 3rd ed. Englewood Cliffs, New Jersey, Prentice-Hall, Inc. 1967.
8. Grant, Eugene L. and Norton, Paul T., Jr. Depreciation. Revised edition. New York, New York, The Ronald Press Co. 1955.
9. Graybill, Franklin A. An introduction to linear statistical models. Vol. 1. New York, New York, McGraw-Hill Book Company, Inc. 1961.
10. Griffen, Daniel L., Jr. Engineering valuation decisions in Iowa and other states. Unpublished M.S. thesis. Ames, Iowa, Library, Iowa State University of Science and Technology. 1961.
11. Guest, P. G. Numerical methods of curve fitting. Cambridge, England, The Syndics of the Cambridge University Press. 1961.
12. Hastings, Cecil, Jr. Approximations for digital computers. Princeton, New Jersey, Princeton University Press. 1955.
13. Henderson, Allen James. The Weibull distribution and industrial property mortality experience. Unpublished M.S. thesis. Ames, Iowa, Library, Iowa State University of Science and Technology. 1965.

14. Hoover, Harold Monroe, Jr. Industrial property life analysis with an analog computer. Unpublished M.S. thesis. Ames, Iowa, Library, Iowa State University of Science and Technology. 1967.
15. Howard, Russell Lewis. Depreciation of railway freight cars. Unpublished M.S. thesis. Ames, Iowa, Library, Iowa State University of Science and Technology. 1947.
16. Johnson, Norman L. and Leone, Fred C. Statistics and experimental design. Vol. 1. New York, New York, John Wiley and Sons, Inc. 1964.
17. Kimball, Bradford F. A system of life tables for physical property based on the truncated normal distribution. *Econometrica* 15: 342-360. 1947.
18. Krane, Scott A. Analysis of survival data by regression techniques. *Technometrics* 5: 161-174. 1963.
19. Kurtz, Edwin B. Life expectancy of physical property. New York, New York, The Ronald Press Company. 1930.
20. Marston, Anson, Winfrey, Robley, and Hempstead, Jean C. Engineering valuation and depreciation. 2nd ed. New York, New York, McGraw-Hill Book Co., Inc. 1953.
21. National Association of Railroad and Utility Commissioners. Report of committee on depreciation. Washington, D.C., National Association of Railroad and Utility Commissioners. 1943.
22. National Association of Railroad and Utility Commissioners. Report of special committee on depreciation. New York, New York, The State Law Reporting Co. 1938.
23. Nichols, Richard Lee. The moment-ratio method of analyzing industrial property experience. Unpublished M.S. thesis. Ames, Iowa, Library, Iowa State University of Science and Technology. 1961.
24. Ostle, Bernard. Statistics in research. Ames, Iowa, The Iowa State University Press. 1954.
25. Sampford, M. R. An introduction to sampling theory. Edinburgh, England, Oliver and Boyd. 1962.
26. Scigliano, J. Michael. An evaluation of the Weibull hazard function as an estimator of industrial property mortality. Unpublished M.S. thesis. Ames, Iowa, Library, Iowa State University of Science and Technology. 1965.

27. Shapiro, S. S. and Wilk, M. B. An analysis of variance test for normality (complete samples). *Biometrika* 52, 591-611. 1965.
28. Winfrey, Robley. Statistical analysis of industrial property retirements: revised April, 1967 by Harold A. Cowles, Professor, Department of Industrial Engineering. Ames, Iowa, Iowa State University of Science and Technology, Engineering Research Institute Bulletin 125, revised edition. 1967.
29. Winfrey, Robley, and Kurtz, E. B. Life characteristics of physical property. Ames, Iowa, Iowa State University of Science and Technology, Engineering Research Institute Bulletin 103. 1931.

ACKNOWLEDGMENTS

The author is deeply indebted to Dr. Harold A. Cowles and Dr. Herbert T. David for their advice, counsel, and encouragement in developing the ideas presented herein. Thanks are also due to Dr. Edward Carney for assistance with the computer programs. The author would also like to acknowledge the patience and understanding of his wife, Anne, and his daughter, Jill-bette, during this period of graduate study.

APPENDIX A - GENERAL FLOW CHART OF SIMULATION PROGRAM

Simulation of the retirement experience of a given vintage group can be used to calculate one retirement ratio for each age interval. Repeated simulation of the retirement experience of the same vintage group yields addition retirement ratios for each age interval. These retirement ratios for an age interval constitute an empirical, vertical distribution of retirement ratios for that age interval. A general flow chart of the computer program to accomplish this simulation of the vertical distribution of retirement ratios at each age interval is presented in this appendix.

An Iowa type curve was used to provide the parent population of ages of units at retirement for the purpose of simulating the retirement experience of a vintage group.

The abbreviations used in the flow chart are:

Arr. = arrange

Betw. = between

Calc. = calculate

Corresp. = corresponding

Cum. = cumulative

Distr. = distribution

Exp. = experience

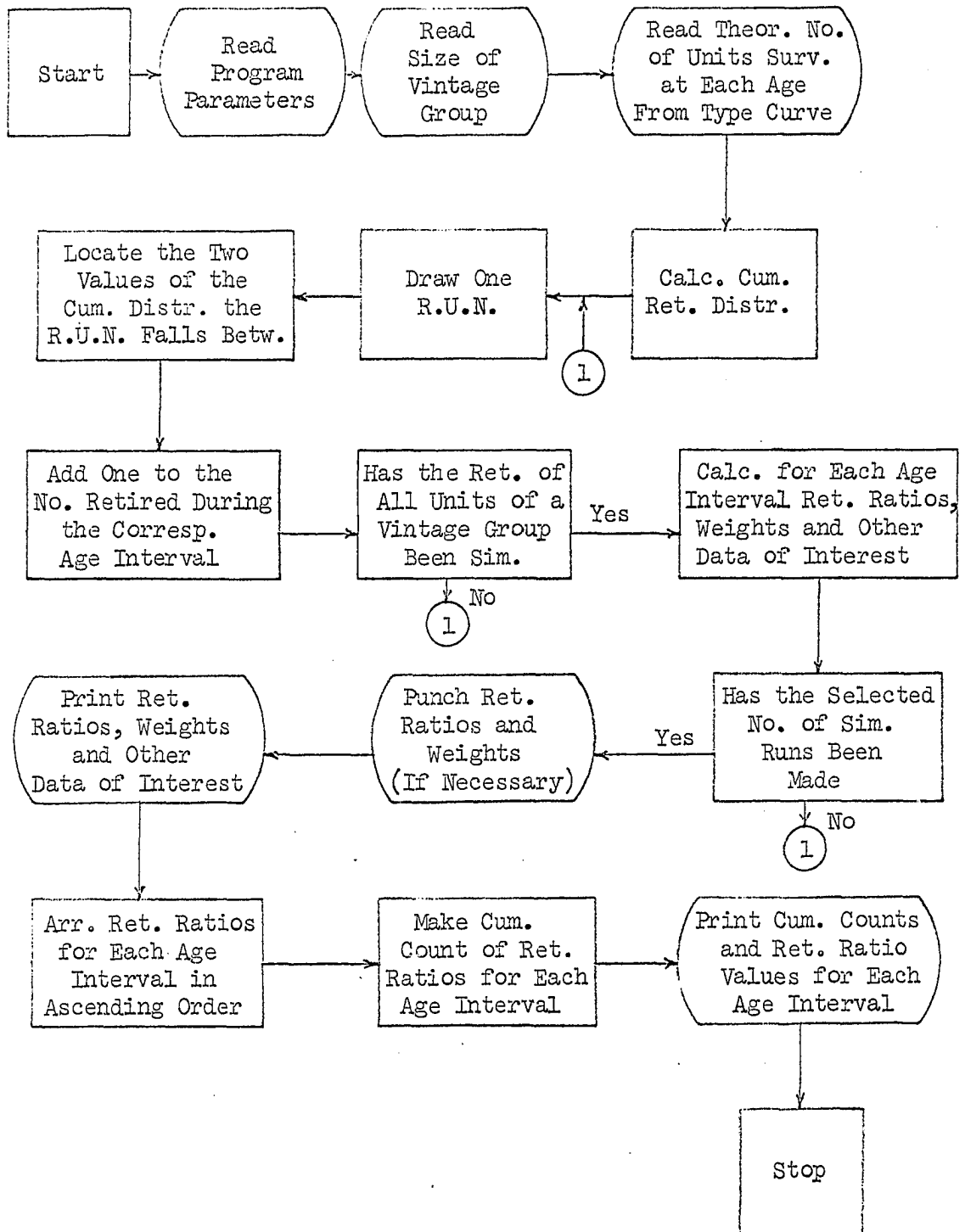
No. = number

Ret. = retirement

R.U.N. = random number from a uniform distribution

Sim. = simulated

Theor. = theoretical



APPENDIX B - GENERAL FLOW CHART OF NORMAL APPROXIMATION PROGRAM

The computer program, the general flow chart of which is presented in this section, generates an approximation of the vertical, cumulative distribution of retirement ratios at each age interval. A vital section of the program is the subroutine, obtained from the Iowa State University Statistical Laboratory - Numerical Analysis and Programming Section, for computing an approximation of the normal cumulative distribution. This subroutine is based on the work of Hastings (12, p. 168).

The abbreviations and symbols used in the flow chart are:

Approx. = approximate

Calc. = calculate

Cum. = cumulative

Distr. = distribution

Incr. = increments

Neg. = negative

No. = number

PR. = $\Pr[(1 - T)L_{\cdot k} - T M_{\cdot k} \leq 0]$

Ret. = retirement

Subr. = subroutine

Surv. = surviving

Theor. = theoretical

$$L_{\cdot k} = \sum_{j=1}^J L_{jk}$$

$L_{jk} = 1$ if the j^{th} unit of the sample is retired during the k^{th} age interval

= 0 otherwise

$$M_{\cdot k} = \sum_{j=1}^J M_{jk}$$

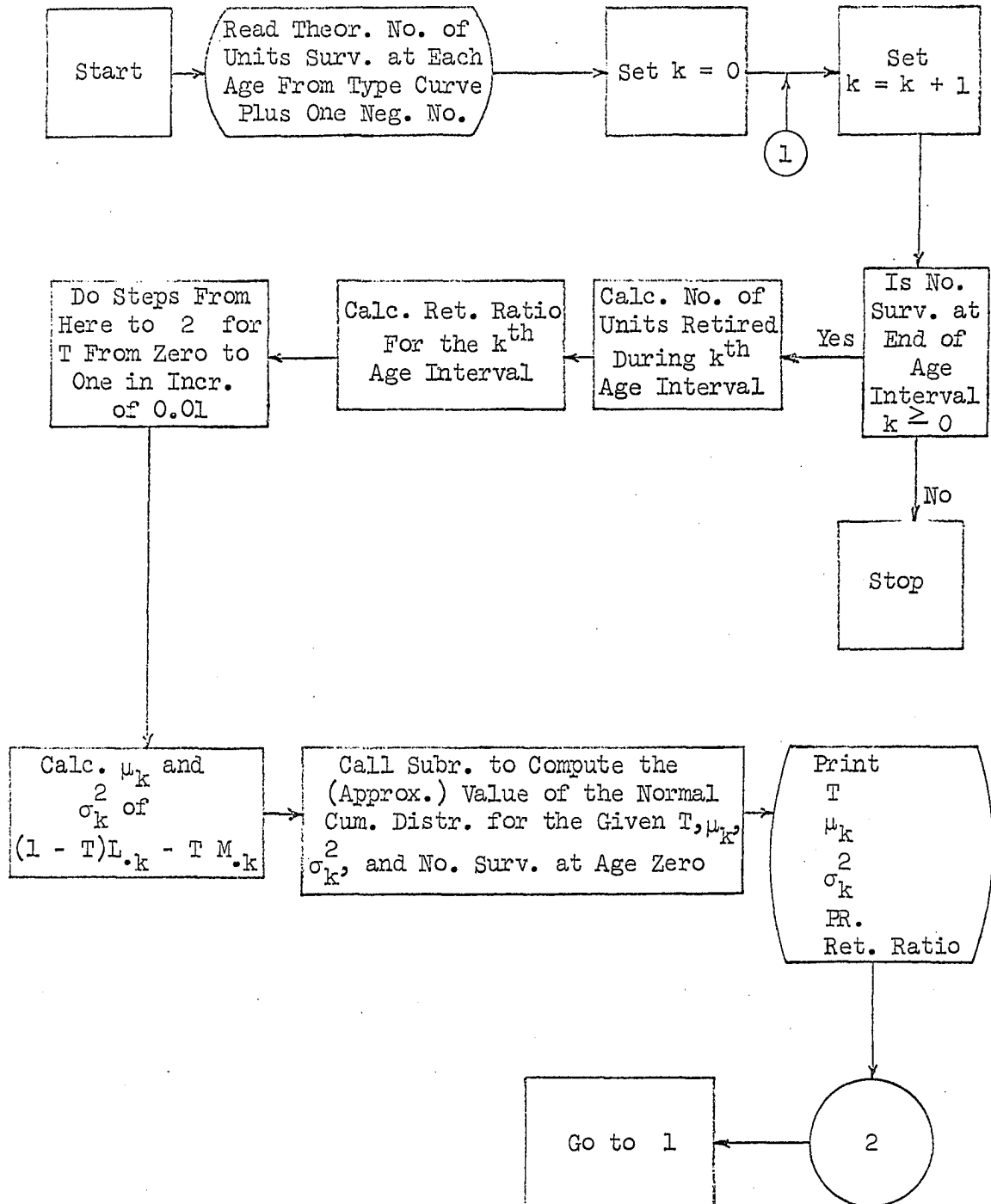
$M_{jk} = 1$ if the j^{th} unit of the sample is retired after the k^{th} age interval

= 0 otherwise

μ_k = mean

σ_k^2 = variance

T = dummy variable



APPENDIX C - MAXIMUM-LIKELIHOOD ESTIMATORS
OF THE POLYNOMIAL COEFFICIENTS

The procedure for determining a maximum-likelihood estimator is

(1, p. 101):

1. Determine the distribution function of the sample,
 $f(X_1, X_2, \dots, X_n; \theta)$.
2. Determine $L = \log f(X_1, X_2, \dots, X_n; \theta)$.
3. Determine the value of θ which will maximize L by solving the equation: $\partial L / \partial \theta = 0$. This will also maximize the likelihood.

Let

μ_k = population retirement ratio for the k^{th} age interval

r_{ik} = retirement ratio for the k^{th} age interval obtained from
the i^{th} vintage group (assumed to be $NID(\mu_k, \sigma_{ik})$)

σ_{ik}^2 = variance of the population of retirement ratios for the
 k^{th} age interval from samples of the size of vintage
group i

$\tilde{\mu}_k$ = maximum-likelihood estimator of μ_k

$\sigma_{\cdot k}^2$ = variance of $\tilde{\mu}_k$.

Then

$$f(r_{ik}) = \frac{1}{\sigma_{ik} \sqrt{2\pi}} e^{-\frac{(r_{ik} - \mu_k)^2}{2\sigma_{ik}^2}}$$

if the r_{ik} are assumed to be normally distributed.

$$r_{ik} \sim N(\mu_k, \sigma_{ik}) \quad i = 1, 2, \dots, I$$

$$k = 1, 2, \dots, K$$

The distribution function of the sample retirement ratios at age interval k is

$$f(r_{1k}, r_{2k}, \dots, r_{Ik}; \mu_k) \\ = \left(\frac{1}{\sqrt{2\pi}}\right)^I \left(\frac{1}{\pi} \frac{1}{\sigma_{ik}}\right) e^{-1/2 \sum_{i=1}^I \left(\frac{r_{ik} - \mu_k}{\sigma_{ik}}\right)^2}$$

where

$$\frac{1}{\pi} \frac{1}{\sigma_{ik}} = \left(\frac{1}{\sigma_{1k}}\right) \left(\frac{1}{\sigma_{2k}}\right) \dots \left(\frac{1}{\sigma_{Ik}}\right)$$

L is the natural log of the distribution function.

$$L = \ln[f(r_{1k}, r_{2k}, \dots, r_{Ik}; \mu_k)] \\ = I \ln\left(\frac{1}{\sqrt{2\pi}}\right) + \ln\left(\frac{1}{\pi} \frac{1}{\sigma_{ik}}\right) - 1/2 \sum_{i=1}^I \left(\frac{r_{ik} - \mu_k}{\sigma_{ik}}\right)^2$$

where $\prod_{i=1}^I$ and $\sum_{i=1}^I$ denote $\prod_{i=1}^I$ and $\sum_{i=1}^I$, respectively. L is maximized by setting the partial derivative of L with respect to μ_k equal to zero and solving for μ_k .

$$\frac{\partial L}{\partial \mu_k} = -1/2(2) \sum_{i=1}^I \left(\frac{r_{ik} - \mu_k}{\sigma_{ik}}\right) \left(-\frac{1}{\sigma_{ik}}\right) = 0$$

$$\sum_{i=1}^I \left(\frac{r_{ik} - \mu_k}{2\sigma_{ik}}\right) = 0$$

$$\sum_{i=1}^I \left(\frac{r_{ik}}{2\sigma_{ik}}\right) = \sum_{i=1}^I \left(\frac{\mu_k}{2\sigma_{ik}}\right)$$

Since μ_k is a constant over all I , by assumption, then

$$\mu_k = \frac{\sum_{ik}^I r_{ik} / \sigma_{ik}^2}{\sum_{ik}^I 1 / \sigma_{ik}^2}$$

$$= \frac{\sum_{ik}^I (1 / \sigma_{ik}^2) (r_{ik})}{\sum_{ik}^I (1 / \sigma_{ik}^2)}$$

Thus, the maximum-likelihood estimator of the retirement ratio for age interval k is the weighted average retirement ratio. The weight to be given each retirement ratio, r_{ik} , is the inverse of the variance of the retirement ratio, σ_{ik}^2 .

Let

$$w_{ik} = \frac{1}{\sigma_{ik}^2}$$

$$w_{\cdot k} = \sum_{ik}^I w_{ik}$$

$$= \sum_{ik}^I (1 / \sigma_{ik}^2)$$

The estimator of the variance of $\tilde{\mu}_k$ is

$$\sigma_{\cdot k}^2 = \text{var}(\tilde{\mu}_k)$$

$$= \text{var}\left(\frac{\sum_{ik}^I w_{ik} r_{ik}}{w_{\cdot k}}\right)$$

$$= \text{var}\left(\frac{w_{1k} r_{1k}}{w_{\cdot k}} + \frac{w_{2k} r_{2k}}{w_{\cdot k}} + \dots + \frac{w_{Ik} r_{Ik}}{w_{\cdot k}}\right)$$

$$= \text{var}\left(\frac{w_{1k} r_{1k}}{w_{\cdot k}}\right) + \text{var}\left(\frac{w_{2k} r_{2k}}{w_{\cdot k}}\right) + \dots + \text{var}\left(\frac{w_{Ik} r_{Ik}}{w_{\cdot k}}\right)$$

For a given sample retirement ratio, r_{ik} , from an $N(\mu_k, \sigma_{ik})$ distribution, the variance of the sample, σ_{ik}^2 , is a constant and, therefore, w_{ik} and $w_{\cdot k}$

are constants. Also, the variance of a constant times a variable is the square of the constant times the variance of the variable. Therefore

$$\sigma_{\cdot k}^2 = \frac{w_{1k}^2}{w_{\cdot k}^2} \text{var}(r_{1k}) + \frac{w_{2k}^2}{w_{\cdot k}^2} \text{var}(r_{2k}) + \dots + \frac{w_{Ik}^2}{w_{\cdot k}^2} \text{var}(r_{Ik})$$

But

$$\begin{aligned} \text{var}(r_{ik}) &= \sigma_{ik}^2 \\ &= \frac{1}{w_{ik}} \end{aligned}$$

Therefore

$$\begin{aligned} \sigma_{\cdot k}^2 &= \left(\frac{w_{1k}^2}{w_{\cdot k}^2}\right)\left(\frac{1}{w_{1k}}\right) + \left(\frac{w_{2k}^2}{w_{\cdot k}^2}\right)\left(\frac{1}{w_{2k}}\right) + \dots + \left(\frac{w_{Ik}^2}{w_{\cdot k}^2}\right)\left(\frac{1}{w_{Ik}}\right) \\ &= \frac{w_{1k}}{w_{\cdot k}^2} + \frac{w_{2k}}{w_{\cdot k}^2} + \dots + \frac{w_{Ik}}{w_{\cdot k}^2} \\ &= \sum \frac{I w_{ik}}{w_{\cdot k}^2} \\ &= \frac{w_{\cdot k}}{w_{\cdot k}^2} \\ &= 1/w_{\cdot k} \\ &= \frac{1}{\sum w_{ik}} \\ &= \frac{1}{\sum (1/\sigma_{ik}^2)} \end{aligned}$$

Thus, the maximum-likelihood estimators of the μ_k , $\tilde{\mu}_k$, and the variances of the $\tilde{\mu}_k$, $\sigma_{\cdot k}^2$, are

$$\tilde{\mu}_k = \frac{\sum (1/\sigma_{ik}^2)(r_{ik})}{\sum (1/\sigma_{ik}^2)}$$

$$\sigma_{\cdot k}^2 = \frac{1}{\sum (1/\sigma_{ik}^2)}$$

The regression equation of retirement ratios on age intervals is, for the first degree case

$$\mu_k = \alpha + \beta x_k + \epsilon_k$$

$$r_{\cdot k} = a + bx_k + e_k$$

where

$r_{\cdot k}$ = "observed" weighted average retirement ratio for age interval k (i.e., the sample calculation of $\tilde{\mu}_k$)

α, β = regression coefficients

ϵ_k = error term

a, b, e_k = estimators of $\alpha, \beta, \epsilon_k$, respectively

x_k = age interval index number

Since $\tilde{\mu}_k$ is distributed $NID(\mu_k, \sigma_{\cdot k})$ (11, pp. 29-30) and, therefore,

$$\epsilon_k \sim NID(0, \sigma_{\epsilon_k})$$

the maximum-likelihood estimators of α and β , a and b , respectively, can be determined. The distribution function of the e_k 's is

$$\begin{aligned} & f(e_1, e_2, \dots, e_K; \epsilon_k) \\ &= \left(\frac{1}{\sqrt{2\pi}}\right)^K \left(\prod \sigma_{\epsilon_k}\right) e^{-1/2 \sum \left(\frac{e_k - \mu_{\epsilon_k}}{\sigma_{\epsilon_k}}\right)^2} \end{aligned}$$

Since the mean of ϵ_k equals zero

$$\begin{aligned}
 & f(e_{.1}, \dots, e_K; \epsilon_K) \\
 & \quad - 1/2 \sum_{\epsilon_K}^K \left(\frac{\epsilon_K}{\sigma} \right)^2 \\
 & = \left(\frac{1}{\sqrt{2\pi}} \right)^K (\pi \sigma_{\epsilon_K})^K e
 \end{aligned}$$

The regression equation, in terms of ϵ_K , is

$$\epsilon_K = \mu_K - \alpha - \beta x_K$$

Therefore

$$\begin{aligned}
 & f(r_{.1}, \dots, r_{.K}, x_1, \dots, x_K; \alpha, \beta) \\
 & \quad - 1/2 \sum_{\epsilon_K}^K \left(\frac{\mu_K - \alpha - \beta x_K}{\sigma} \right)^2 \\
 & = \left(\frac{1}{\sqrt{2\pi}} \right)^K (\pi \sigma_{\epsilon_K})^K e
 \end{aligned}$$

The L is

$$\begin{aligned}
 L & = \ln[f(r_{.1}, \dots, r_{.K}, x_1, \dots, x_K; \alpha, \beta)] \\
 & = K \ln \left(\frac{1}{\sqrt{2\pi}} \right) + \ln(\pi \sigma_{\epsilon_K})^K - 1/2 \sum_{\epsilon_K}^K \left(\frac{\mu_K - \alpha - \beta x_K}{\sigma} \right)^2
 \end{aligned}$$

The maximum-likelihood estimator of α is given by

$$\frac{\partial L}{\partial \alpha} = -1/2 \sum_{\epsilon_K}^K \left(\frac{\mu_K - \alpha - \beta x_K}{\sigma} \right) \left(-\frac{1}{\sigma} \right) = 0$$

$$\sum_{\epsilon_K}^K \mu_K / \sigma_{\epsilon_K}^2 = \sum_{\epsilon_K}^K \alpha / \sigma_{\epsilon_K}^2 + \sum_{\epsilon_K}^K \beta x_K / \sigma_{\epsilon_K}^2$$

$$\sum_{\epsilon_K}^K \mu_K / \sigma_{\epsilon_K}^2 = \alpha \sum_{\epsilon_K}^K 1 / \sigma_{\epsilon_K}^2 + \beta \sum_{\epsilon_K}^K x_K / \sigma_{\epsilon_K}^2 \quad (3)$$

The maximum-likelihood estimator of β is given by

$$\frac{\partial \bar{L}}{\partial \beta} = -1/2 \sum_{k=1}^K \left(\frac{\mu_k - \alpha - \beta x_k}{\sigma_{\epsilon_k}^2} \right) \left(-\frac{x_k}{\sigma_{\epsilon_k}^2} \right)$$

$$\sum_{k=1}^K \mu_k x_k / \sigma_{\epsilon_k}^2 = \alpha \sum_{k=1}^K x_k / \sigma_{\epsilon_k}^2 + \beta \sum_{k=1}^K x_k^2 / \sigma_{\epsilon_k}^2 \quad (4)$$

Since a and b are constants and x is assumed to be measured without error, the variance of $r_{\cdot k}$, $\sigma_{\cdot k}^2$, is the variance of ϵ_k , $\sigma_{\epsilon_k}^2$. Replacing μ_k , α , β , and $\sigma_{\epsilon_k}^2$ by $r_{\cdot k}$, a , b , and $\sigma_{\cdot k}^2$, respectively, equations 3 and 4 become

$$\sum_{k=1}^K r_{\cdot k} / \sigma_{\cdot k}^2 = a \sum_{k=1}^K 1 / \sigma_{\cdot k}^2 + b \sum_{k=1}^K x_k / \sigma_{\cdot k}^2 \quad (5)$$

$$\sum_{k=1}^K x_k r_{\cdot k} / \sigma_{\cdot k}^2 = a \sum_{k=1}^K x_k / \sigma_{\cdot k}^2 + b \sum_{k=1}^K x_k^2 / \sigma_{\cdot k}^2 \quad (6)$$

or

$$\sum_{k=1}^K w_{\cdot k} r_{\cdot k} = a \sum_{k=1}^K w_{\cdot k} + b \sum_{k=1}^K w_{\cdot k} x_k$$

$$\sum_{k=1}^K x_k r_{\cdot k} w_{\cdot k} = a \sum_{k=1}^K x_k w_{\cdot k} + b \sum_{k=1}^K x_k^2 w_{\cdot k}$$

which are the same as the first two normal equations obtained by the principle of least-squares.

The maximum-likelihood estimators for the coefficients of higher degree polynomials can be solved for in a similar manner.

If the variances are known, the maximum-likelihood estimators of the polynomial coefficients are unbiased and have minimum variance amongst all unbiased estimators (9, pp. 113-114, 117).

APPENDIX D - PRELIMINARY APPROACH

Dr. Fuller¹ suggested a preliminary approach (to the problem of fitting a polynomial to the retirement ratios) based on sampling theory. The link between retirement ratios and sampling theory is the analogy between the several retirement ratios at an age interval (one from each vintage group) and cluster sampling from proportions (25, pp. 236-237).

Let

P_k = population retirement ratio for the k^{th} age interval

r_{ik} = sample estimate of P_k from the i^{th} vintage group (or i^{th} cluster)

i = vintage group (or cluster) index number

= 1, 2, . . . ; I

j = unit number within a vintage group (or cluster)

= 1, 2, . . . , J

k = age interval index number

= 1, 2, . . . , K

L_{ijk} = 1 if the j^{th} item of the i^{th} vintage group (or i^{th} cluster) is retired during the k^{th} age interval

= 0 otherwise

M_{ijk} = 1 if the j^{th} item of the i^{th} vintage group (i^{th} cluster) is surviving at the beginning of the k^{th} age interval

= 0 otherwise

$J_{ik} = \sum_j M_{ijk}$

= $M_{i \cdot k}$

¹ Fuller, Wayne A., Professor of Statistics Department, Iowa State University of Science and Technology, Ames, Iowa. Information on sampling theory. Private communication. 1967.

= number of units from the i^{th} vintage group surviving at the beginning of the k^{th} age interval (or size of the i^{th} cluster at the k^{th} age interval).

$\tilde{P}_{\cdot k}$ = composite estimate of P_k

Then

$$\begin{aligned} r_{ik} &= \frac{\sum_{j=1}^J L_{ijk}}{\sum_{j=1}^J M_{ijk}} \\ &= \frac{\sum_{j=1}^J L_{ijk}}{J_{ik}} \end{aligned}$$

A single estimate, $\tilde{r}_{\cdot k}$, of the population retirement ratio for age interval k can be obtained by weighting each r_{ik} by the inverse of its variance. The variance of r_{ik} , conditional upon the denominator of r_{ik} , is

$$\begin{aligned} \text{var}(r_{ik}) &= \text{var}\left(\frac{\sum_{j=1}^J L_{ijk}}{J_{ik}}\right) \\ &= \frac{1}{J_{ik}^2} \text{var}\left(\sum_{j=1}^J L_{ijk}\right) \\ &= \frac{1}{J_{ik}^2} J_{ik} P_k Q_k \end{aligned}$$

since L_{ijk} is binomially distributed. Therefore

$$\text{var}(r_{ik}) = \frac{P_k Q_k}{J_{ik}}$$

and

$$w_{ik} = \frac{1}{\text{var}(r_{ik})}$$

$$= \frac{J_{ik}}{P_k Q_k}$$

Then

$$\begin{aligned} \tilde{r}_{\cdot k} &= \frac{\sum_{ik} w_{ik} r_{ik}}{\sum_{ik} w_{ik}} \\ &= \frac{\sum_{ik} \frac{J_{ik}}{P_k Q_k} r_{ik}}{\sum_{ik} \frac{J_{ik}}{P_k Q_k}} \\ &= \frac{\sum_{ik} J_{ik} r_{ik}}{\sum_{ik} J_{ik}} \end{aligned}$$

since P_k and Q_k are constants over all i for a given k .

The weight to give each $\tilde{r}_{\cdot k}$ when fitting a polynomial to the $\tilde{r}_{\cdot k}$ is not clear. A suggested weight for each $\tilde{r}_{\cdot k}$ is

$$\tilde{w}_{\cdot k} = \frac{\sum_{ik} J_{ik}}{P_k Q_k}$$

where

1. $\tilde{w}_{\cdot k}$ is the inverse of the variance of $\tilde{r}_{\cdot k}$,
2. The variance in (1) is conditional upon the denominators, S_{ik} , and
3. $\tilde{r}_{\cdot k}$ is the composite sample estimate of P_k .

The expression to minimize in fitting a polynomial to the retirement ratios by a weighted least squares approach is, then,

$$\text{Min}_{a,b,c,\text{etc.}} \left\{ \sum_{ik} \frac{J_{ik}}{\tilde{r}_{\cdot k} (1 - \tilde{r}_{\cdot k})} [\tilde{r}_{\cdot k} - (a + bx_k + cx_k^2 + \dots)]^2 \right\}$$

where

$\tilde{r}_{.k}$ = sample estimate of P_k

The variances of the $\tilde{r}_{.k}$ used in fitting a polynomial to the retirement ratios, $\tilde{\sigma}_{.k}^2$, present a theoretical problem. Since

1. Each variance is conditional upon the denominator,
2. The numerators and denominators are related in the following manner

$$\tilde{r}_{.1} = \frac{L_{..1}}{M_{..1}}$$

$$\tilde{r}_{.2} = \frac{L_{..2}}{M_{..2}}$$

$$M_{..2} = M_{..1} - L_{..1}$$

$$L_{..1} = M_{..1} - M_{..2}$$

$$\tilde{r}_{.3} = \frac{L_{..3}}{M_{..3}}$$

$$M_{..3} = M_{..2} - L_{..2}$$

$$L_{..2} = M_{..2} - M_{..3}$$

.

.

.

$$\tilde{r}_{.K} = \frac{L_{..K}}{M_{..K}}$$

$$M_{..K} = M_{..K-1} - L_{..K-1}$$

$$L_{..K} = M_{..K}$$

$$\sum_{IJK} \sum_{\Sigma} L_{ijk} = L_{...}$$

$$= \sum_{IJ} \sum_{\Sigma} M_{ij1}$$

$$= M_{..1}$$

then the variances, when all J_{ik} are known and considered, are variances of constants and are zero. The extent to which this theoretical consideration limits the usefulness of the practical application of the procedure is not known.

The least squares expression of this preliminary approach

$$\text{Min}_{a,b,c,\text{etc.}} \left\{ \sum_{k=1}^K \frac{J_{\cdot k}}{\bar{r}_{\cdot k}(1 - \bar{r}_{\cdot k})} [\bar{r}_{\cdot k} - (a + bx_k + cx_k^2 + \dots)]^2 \right\}$$

is quite similar to the least squares expression of a present method of obtaining a smoothed life table (see p. 54 of this dissertation)

$$\text{Min}_{a,b,c,\text{etc.}} \sum_{k=1}^K \{ S_{\cdot k} [\bar{r}_{\cdot k} - (a + bx_k + cx_k^2 + \dots)]^2 \}$$

where the $\bar{r}_{\cdot k}$ in both expressions is

$$\bar{r}_{\cdot k} = \frac{\sum_{i=1}^I \text{retirements during the } k^{\text{th}} \text{ age interval from vintage group } i}{\sum_{i=1}^I \text{survivors at the beginning of the } k^{\text{th}} \text{ age interval from vintage group } i}$$

and

$$S_{\cdot k} = J_{\cdot k} = \sum_{i=1}^I \text{survivors at the beginning of the } k^{\text{th}} \text{ age interval from vintage group } i$$

The import of this similarity between the least squares expressions would seem to be that the present method of fitting a polynomial to the weighted average retirement ratio

$$\bar{r}_{\cdot k} = \frac{\sum_{i=1}^I R_{ik}}{\sum_{i=1}^I S_{ik}}$$

by a least-squares approach where each $\tilde{r}_{.k}$ is weighted by

$$\tilde{w}_{.k} = S_{.k}$$

may be a reasonable method, but a method in which the weight given each $\tilde{r}_{.k}$ could, perhaps, be improved upon.

APPENDIX E - GENERAL FLOW CHART OF PROGRAM TO IMPLEMENT THE PROCEDURE

The computer program to implement the polynomial fitting procedure developed in this dissertation is actually a combination of several computer programs. The basic parts of the program are the subroutine to fit polynomials to retirement ratios, the subroutine to compute an approximation of the normal, cumulative distribution, and the section which computes σ_{ik}^2 . The remaining parts of the program primarily process the data to obtain the necessary input to the above-mentioned subroutines and provide instructions to the computer as to when to proceed to which operations.

The abbreviations and symbols used in the flow chart are:

Approx. = approximate

Calc. = calculate

Cum. = cumulative

Deg. = degree(s)

Distr. = distribution

Extrap. = extrapolate

Inc. = increment(s)

Interp. = interpolate

No. = number

Polyn. = polynomial(s)

Ret. = retirement

Sign. = significance, significant

Subr. = subroutine

Surv. = surviving

Wtg. = weighting

$$\tilde{r}_{\cdot k} = \frac{L_{\cdot \cdot k}}{L_{\cdot \cdot k} + M_{\cdot \cdot k}}$$

$$L_{\cdot \cdot k} = \sum_{i=1}^I \sum_{j=1}^J L_{ijk}$$

$L_{ijk} = 1$ if the j^{th} unit of the i^{th} vintage group is retired during the k^{th} age interval
 $= 0$ otherwise

$$M_{\cdot \cdot k} = \sum_{i=1}^I \sum_{j=1}^J M_{ijk}$$

$M_{ijk} = 1$ if the j^{th} unit of the i^{th} vintage group is retired after the k^{th} age interval
 $= 0$ otherwise

$$\tilde{w}_{\cdot k} = \frac{L_{\cdot \cdot k} + M_{\cdot \cdot k}}{\tilde{r}_{\cdot k}(1 - \tilde{r}_{\cdot k})}$$

$\tilde{r}_{\cdot k}$ = value of the retirement ratio for the k^{th} age interval interpolated from the polynomial fit of the $\tilde{r}_{\cdot k}$

T = dummy variable

Δ = delta; amount of increment

$$C_k = \Pr(L_{i \cdot k} = 1)$$

$$C'_k = \Pr(M_{i \cdot k} = 1)$$

μ_k = mean

σ_k^2 = variance

$$P = \Pr[(1 - T)L_{i \cdot k} - T L_{i \cdot k} \leq 0]$$

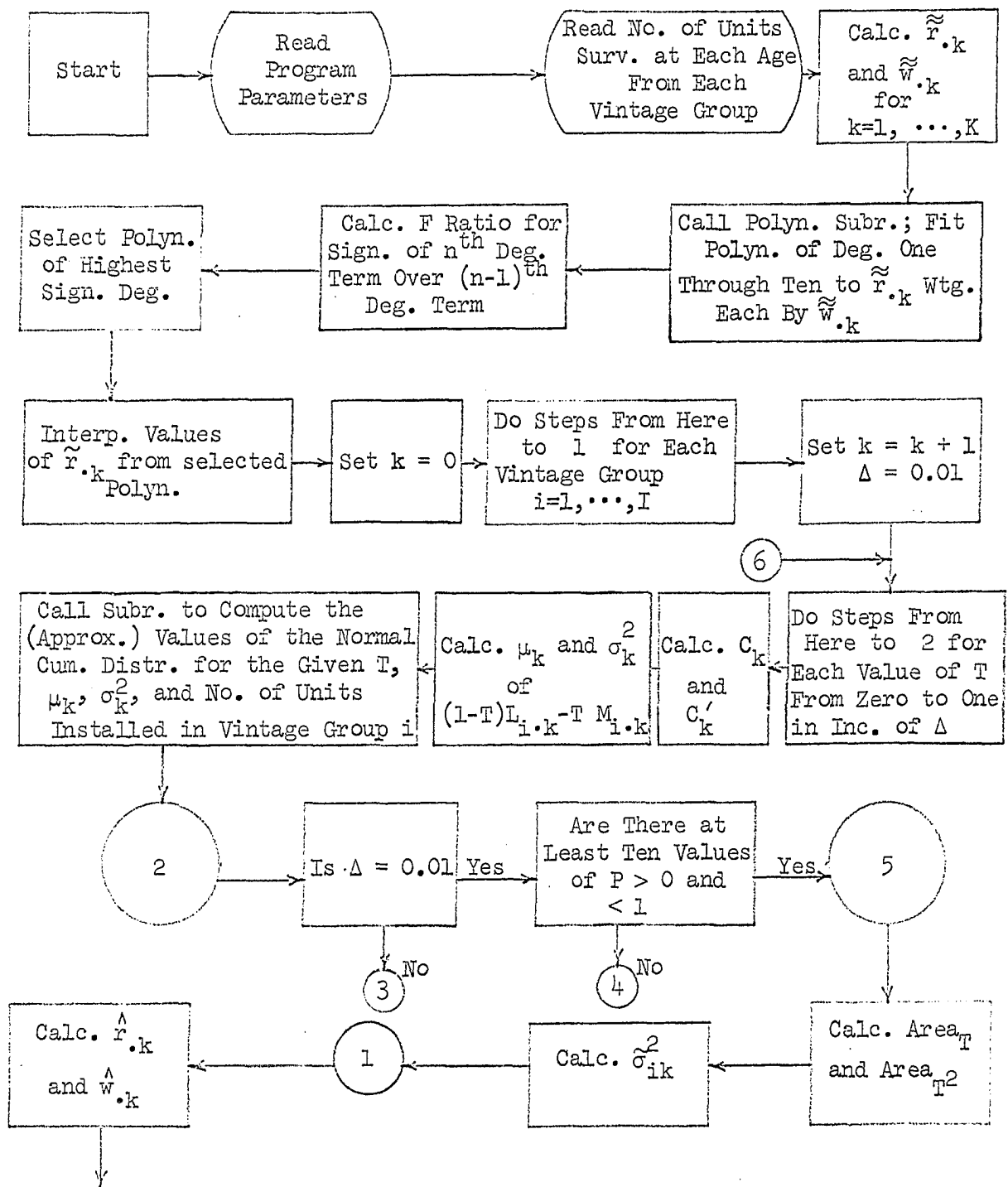
Area_T = area above the cumulative distribution versus
T curve

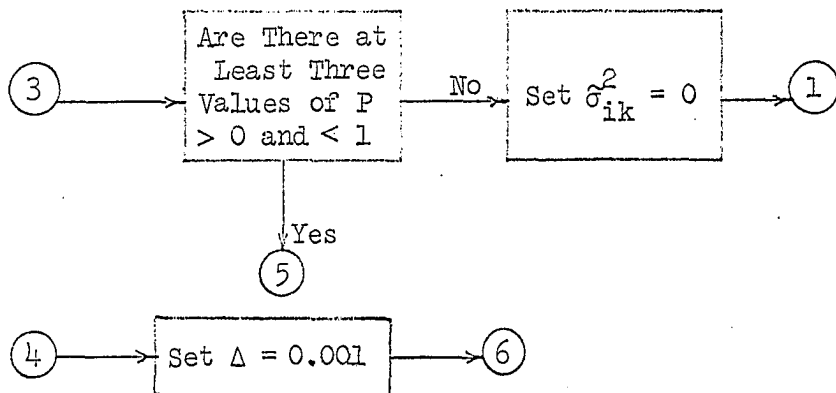
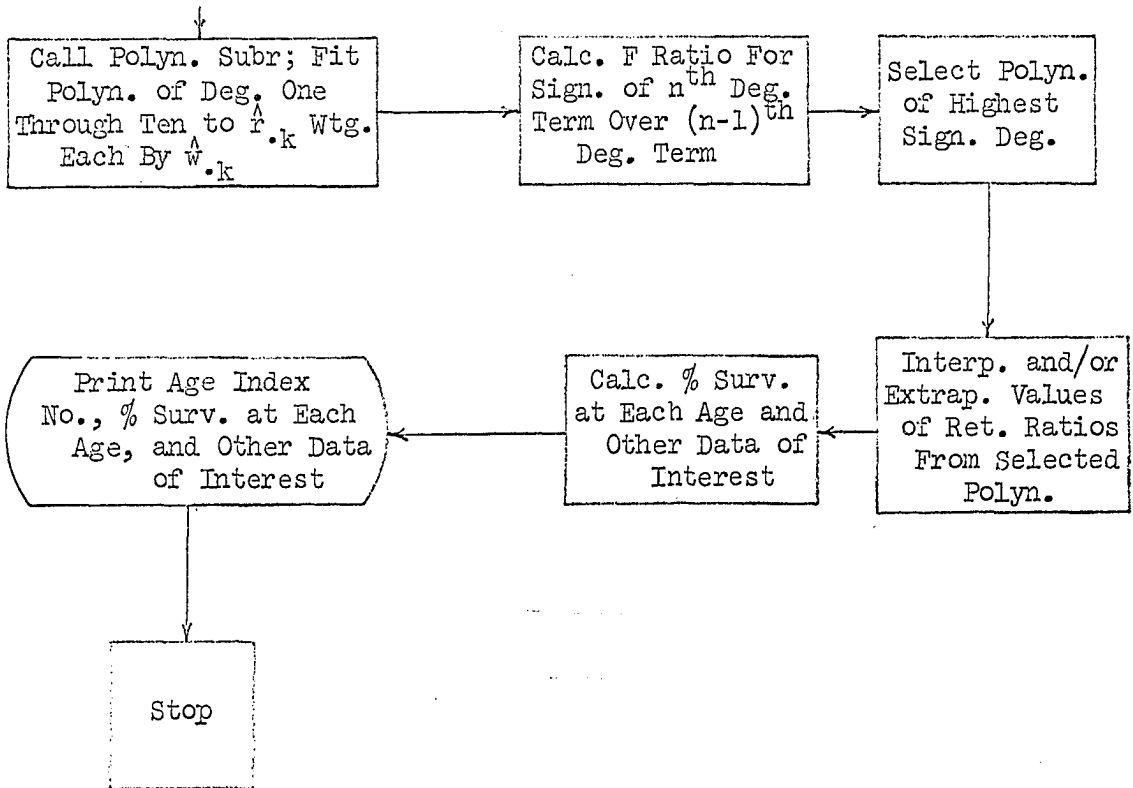
Area_{T²} = area above the cumulative distribution versus
T² curve

$$\hat{r}_{\bullet k} = \frac{\sum_{i=1}^I w_{ik} r_{ik}}{\sum_{i=1}^I w_{ik}}$$

$$\hat{w}_{\bullet k} = \sum_{i=1}^I \frac{1}{\hat{\sigma}_{ik}^2}$$

$$w_{ik} = \frac{1}{\sigma_{ik}^2}$$





APPENDIX F - TESTING FOR NORMALITY

The use of certain statistical procedures, such as setting confidence limits or making tests of significance, requires an assumption about the distribution of the variable. Frequently, the variable is assumed to be normally distributed. Therefore, a considerable amount of research has been done on the problem of testing the normality of a sample (27, p. 591):

Testing for distributional assumptions in general and for normality in particular has been a major area of continuing statistical research--both theoretically and practically. A possible cause of such sustained interest is that many statistical procedures have been derived based on particular distributional assumptions--especially that of normality.

A number of statistics are available for testing the hypothesis

$$H_0: \{Y_i\} \sim N(\mu, \sigma) \quad i = 1, 2, \dots, I$$

where $\{Y_i\}$ is a sample of size I . In testing a hypothesis, two types of error are possible (24, p. 27):

Type I error. If we reject our hypothesis when it is actually true, we have committed an error of the first kind, or a Type I error.

Type II error. If we accept our hypothesis when it is actually false, we have committed an error of the second kind, or a Type II error.

The probability of a Type I error, α , is represented as

$$\Pr(\text{rejecting } H_0 | H_0 \text{ true}) = \alpha$$

The probability of a Type II error, β , is represented as

$$\Pr(\text{accepting } H_0 | H_0 \text{ false}) = \beta$$

$1 - \beta$ is called the power of the test and may be represented as

$$\Pr(\text{reject } H_0 | H_0 \text{ false}) = 1 - \beta$$

Obviously, a test of a H_0 which minimizes both α and β would be desirable, however (24, p. 27):

We shall remark here that, if our size of sample (number of sample observations) has been decided in advance, it is not possible to minimize α and β simultaneously.

A common procedure is to select an α , a sample size, and a test. A desirable attribute of such a test is that (for any given sample size and α) β is equal to or less than the β of any other test (or the power of the test is equal to or greater than the power of any other test). If a test with this attribute is not available, then the test which has optimum over-all power, according to some criterion, should be selected. The simplicity of the test, from an applications point of view, may be an important, additional criterion.

Shapiro and Wilk (27) developed the W test for testing the hypothesis

$$H_0: \{y_i\} \sim N(\mu, \sigma) \quad i = 1, 2, \dots, I$$

and empirically obtained the power of the W and eight other tests against each of fifteen different distributions. Comparing the power of the W statistic with the power of any one of the other eight tests shows that the power of the W statistic is greater against at least a majority of the fifteen distributions.

The W statistic is (27, p. 602-603)

$$W = \frac{b^2}{s^2}$$

where

$$b = \sum_{i=1}^k a_{n-i+1} (y_{n-i+1} - y_i)$$

a = a set of multipliers obtained from a table and
dependent upon n

$$S^2 = \sum_{i=1}^n (y_i - \bar{y})^2$$

y_i = observed values arranged in ascending order

$i = 1, 2, \dots, n$

$$k = \frac{n}{2} \text{ if } n \text{ is even}$$

$$= \frac{n-1}{2} \text{ if } n \text{ is odd.}$$

The W test is origin and scale invariant (27, p. 593).

One of the criteria occasionally used in selecting a test is the simplicity in application of the test. The computation of the W statistic requires a set of "a" factors which are different for different sample sizes. Thus, a table of "a" factors must be available when the W statistic is used.

The objective of the investigation reported in this appendix was to find, by simulation, a test (or tests) for normality simpler than the W test yet having power at least comparable to that of the W test against other distributions.

A possible statistic for testing for normality might be the coefficient of correlation

$$r = \pm \sqrt{1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}} \quad i = 1, 2, \dots, I$$

where

$\sum_{i=1}^I (y_i - \hat{y}_i)^2$ = the sum of the squares of the vertical
deviations from a linear regression line fitted
by the method of least-squares

$\sum_{i=1}^I (y_i - \bar{y})^2$ = the sum of the squares of the deviations of
the observed values from the mean of the
observed values

The statistic (the coefficient of correlation)--is, undoubtedly,
the most widely used measure of the strength of the linear
relationship between two variables. (7, p. 355)

The denominator, $\sum (y_i - \bar{y})^2$, is a measure of the total variation of the
y's. The numerator, $\sum (y_i - \hat{y}_i)^2$, is a measure of the chance variation
(i.e., a measure of the variation not explained by a linear relationship
between x and y).

If a coefficient of correlation type statistic is to be used to test
for normality, selection of the x_i variables in order to obtain a linear
relationship between the y_i and the x_i is crucial. The scheme devised to
obtain a linear relationship between the y_i and the x_i is analogous to
plotting the y_i on normal probability paper. If a random sample is drawn
from a normal population and the sample values arranged in ascending order
and plotted on normal probability paper, the sample values will fall
closely about a straight line. The ordinates of the probability plot are
the ordered y_i (hereafter denoted as Y_i) and they are plotted on a linear
scale. The abscissa values are the percents of the cumulative distribu-
tion and they are plotted according to distances representing the
standard deviates of the normal distribution. Hence, a set of Y_i should
be linearly related to the ordered set of standard deviates (hereafter
denoted as X_i) if the y_i are random samples from a normal population.

Two modifications of the coefficient of correlation were made in
developing the test statistics. Firstly, only the ratio portion of r was

used because the ratio

$$\frac{\sum_{i=1}^I (Y_i - \hat{Y}_i)^2}{\sum_{i=1}^I (Y_i - \bar{Y})^2}$$

is the proportion of the total variation due to chance (7, p. 359) and, therefore, is a measure of the amount of departure of the $\{Y_i\}$ from normality. Secondly, some statistics of the same form but with greater exponents were investigated because ratios of this form with the deviations raised to the third and fourth power are measure of skewness and kurtosis, respectively.

The test statistics investigated were of the form

$$S = \frac{\left[\sum_{i=1}^I |Y_i - \hat{Y}_i|^g \right]^h}{\left[\sum_{i=1}^I |Y_i - \bar{Y}|^c \right]^d} ; \text{ where } gh = dc$$

$\{Y_i\}$ = sample values drawn from some distribution and arranged in ascending order; $i = 1, 2, \dots, I$

I = sample size

\bar{Y} = mean value of Y

$$\bar{Y} = \frac{\sum_{i=1}^I Y_i}{I}$$

$$\hat{Y}_i = a + bX_i$$

$\{X_i\}$ = the ordered standard deviates

where a and b were obtained by fitting a linear equation to the paired values $\{Y_i, X_i\}$.

The S statistics are scale and origin invariant. The proof is as follows:

1. The normal equations for estimating α and β are

$$(\text{"}\Sigma\text{" means } \sum_{i=1}^I)$$

$$\Sigma Y_i = Ia + b \Sigma X_i$$

$$\Sigma X_i Y_i = a \Sigma X_i + b \Sigma (X_i)^2$$

and the regression equation is

$$\hat{Y}_i = a + bX_i$$

2. Solving for a and b yields

$$a = \frac{\Sigma (X_i)^2 \Sigma Y_i - \Sigma X_i \Sigma X_i Y_i}{I \Sigma (X_i)^2 - (\Sigma X_i)^2}$$

$$b = \frac{\Sigma X_i Y_i - \Sigma X_i \Sigma Y_i / I}{\Sigma (X_i)^2 - (\Sigma X_i)^2 / I}$$

3. Since the X_i are the standard deviates (of the normal distribution) and are symmetrically located about the origin, $\Sigma X_i = 0$.

Therefore,

$$\begin{aligned} a &= \frac{\Sigma (X_i)^2 \Sigma Y_i}{I \Sigma (X_i)^2} \\ &= \frac{\Sigma Y_i}{I} \end{aligned}$$

$$b = \frac{\Sigma X_i Y_i}{\Sigma (X_i)^2}$$

4. Assume a set of Y_i are drawn and transformed as follows

$$Y'_i = K(Y_i - f)$$

Then

$$a' = \frac{\Sigma Y'_i}{I}$$

$$= \frac{\sum K(Y_i - f)}{I}$$

and

$$b' = \frac{\sum X_i Y_i'}{\sum (X_i)^2}$$

$$= \frac{\sum X_i K(Y_i - f)}{\sum (X_i)^2}$$

5. The form of the statistics is

$$S = \frac{(\sum |Y_i' - \bar{Y}_i'|^g)^h}{(\sum |Y_i' - \bar{Y}'|^c)^d}; gh = cd$$

$$= \frac{(\sum |Y_i' - (a' + b'X_i)|^g)^h}{(\sum |Y_i' - \bar{Y}'|^c)^d}$$

Then

$$S = \frac{\{\sum |K(Y_i - f) - [\frac{\sum K(Y_i - f)}{I} + X_i \frac{\sum X_i K(Y_i - f)}{\sum (X_i)^2}]|^g\}^h}{\{\sum |K(Y_i - f) - \frac{\sum K(Y_i - f)}{I}|^c\}^d}$$

$$= \frac{K^{gh} \{\sum |Y_i - f - \frac{\sum Y_i}{I} + f - \frac{X_i \sum X_i Y_i}{\sum (X_i)^2} + \frac{X_i f \sum X_i}{\sum (X_i)^2}|^g\}^h}{K^{cd} \{\sum |Y_i - f - \frac{\sum Y_i}{I} + f|^c\}^d}$$

Collecting terms

$$(\sum |Y_i - \frac{\sum Y_i}{I} - \frac{X_i \sum X_i Y_i}{\sum (X_i)^2}|^g)^h$$

$$S = \frac{(\sum |Y_i - \frac{\sum Y_i}{I}|^c)^d}{(\sum |Y_i - \frac{\sum Y_i}{I}|^c)^d}$$

since $\sum X_i = 0$.

6. Let

$$a = \frac{\sum Y_i}{I}$$

$$b = \frac{\sum X_i Y_i}{\sum (X_i)^2}$$

Then, the numerator becomes

$$\begin{aligned} & (\sum |Y_i - [\frac{Y_i}{I} + X_i \frac{\sum X_i Y_i}{\sum X_i}]|^g)^h \\ &= (\sum |Y_i - [a + bX_i]|^g)^h \\ &= (\sum |Y_i - \hat{Y}_i|^g)^h \end{aligned}$$

and

$$S = \frac{(\sum |Y_i - \hat{Y}_i|^g)^h}{(\sum |Y_i - \bar{Y}|^c)^d}$$

thus completing the proof.

The powers of the S test were obtained empirically by simulation. A preliminary simulation run was conducted as follows:

1. A large number of S type statistics were formulated.
2. Sample size was set equal to twenty and α set equal to 0.05.
3. One-hundred samples (each of size twenty) were drawn by simulation from the normal distribution and from each of fifteen different distributions (including fourteen of those utilized by Shapiro and Wilk; see 27, p. 608).
4. The value of each S statistic and the W statistic (for $\alpha = 0.05$) was determined from the samples from the normal distribution.
5. The power of the test of each S statistic against samples from each non-normal distribution was determined and compared with the

power of the test of the W statistic against the same distributions.

6. On the basis of (5), above, nine of the most promising statistics were selected for further study.

The statistics selected for further investigation were

$$s_2 = \frac{\sum |Y_i - \hat{Y}_i|^2}{\sum |Y_i - \bar{Y}|^2}$$

$$s_{12} = \frac{[\sum |Y_i - \hat{Y}_i|^2]^2}{\sum |Y_i - \bar{Y}|^4}$$

$$s_{25} = \frac{[\sum |Y_i - \hat{Y}_i|^2]^3}{\sum |Y_i - \bar{Y}|^6}$$

$$s_{26} = \frac{[\sum |Y_i - \hat{Y}_i|^2]^4}{\sum |Y_i - \bar{Y}|^8}$$

$$s_{13} = \frac{[\sum |Y_i - \hat{Y}_i|^3]^2}{[\sum |Y_i - \bar{Y}|^4]^{3/2}}$$

$$s_{61} = \frac{[\sum |Y_i - \hat{Y}_i|^3]^2}{\sum |Y_i - \bar{Y}|^6}$$

$$s_{28} = \frac{[\sum |Y_i - \hat{Y}_i|^4]^{3/2}}{\sum |Y_i - \bar{Y}|^6}$$

$$s_{40} = \frac{\sum |Y_i - \hat{Y}_i|^8}{\sum |Y_i - \bar{Y}|^8}$$

$$s_{43} = \frac{\sum |Y_i - \hat{Y}_i|^{12}}{\sum |Y_i - \bar{Y}|^{12}}$$

The distributions utilized and the number of samples (of three different sample sizes) drawn from each of the distributions for the final simulation run are shown in Table 5. The distributions utilized by Shapiro and Wilk were all of those shown in Table 5 down to, and including, the $T(10, 3.1)$. A comparison of the power of the W test (for $\alpha = 0.05$ and $I = 20$) against the non-centralized χ^2 distribution obtained by Shapiro and Wilk, 0.59, with the power of the W test (same α and I) obtained in the investigation, 0.15, indicates that the non-centralized χ^2 distributions used were probably not the same distribution.

Table 5. Sample size and number of samples for the final simulation runs

Distribution	Sample size		
	10	20	50
Normal	2000	2000	2000
$\chi^2_{(1)}$	----	1000	----
$\chi^2_{(2)}$	----	1000	----
$\chi^2_{(4)}$	----	1000	----
$\chi^2_{(10)}$	500	500	----
Non-cent. $\chi^2_{(16)}$	----	500	----
Log normal	----	1000	----
Cauchy	1000	1000	1000
Uniform	1000	1000	1000
Logistic	1000	1000	1000
Beta (2, 1)	----	1000	----
LaPlace	----	1000	----
Poisson	----	1000	----
Binomial	1000	1000	1000
$T(5, 2.4)$	----	1000	----

Table 5 (Continued)

Distribution	Sample size		
	10	20	30
T(10, 3.1)	----	1000	----
Half-normal	1000	1000	1000
Half-Cauchy	----	1000	----
Sum of 3 uniforms	----	1000	----

The powers of the tests included in the final simulation run are presented in Tables 6 through 14. Tables 15 through 23 show the differences between the power of the S tests and the W test. A "+" sign indicates that the power of the S test was greater than the power of the W test by the indicated amount; a "-" sign indicates the opposite. Table 24 shows the sum of differences across all α and I values and the largest positive and negative differences.

Table 6. Empirical power of tests for $\alpha = 0.03$ and I = 10

Distribution	W	S2	S12	S25	S26	S13	S61	S28	S40	S43
Binomial	0.33	0.28	0.46	0.42	0.42	0.43	0.44	0.42	0.34	0.30
Uniform	0.05	0.03	0.10	0.12	0.14	0.07	0.11	0.09	0.09	0.10
Cauchy	0.55	0.56	0.52	0.46	0.42	0.54	0.51	0.52	0.51	0.51
Half-normal	0.13	0.12	0.15	0.12	0.10	0.13	0.13	0.13	0.11	0.10
$\chi^2_{(10)}$	0.08	0.09	0.10	0.07	0.06	0.09	0.08	0.09	0.08	0.08
Logistic	0.07	0.07	0.05	0.04	0.03	0.06	0.05	0.05	0.05	0.05

Table 7. Empirical power of tests for $\alpha = 0.05$ and $I = 10$

Distribution	W	S2	S12	S25	S26	S13	S61	S28	S40	S43
Binomial	0.48	0.39	0.50	0.53	0.52	0.50	0.52	0.51	0.48	0.47
Uniform	0.09	0.05	0.14	0.18	0.19	0.11	0.15	0.13	0.13	0.13
Cauchy	0.58	0.59	0.56	0.52	0.49	0.57	0.55	0.55	0.54	0.53
Half-normal	0.18	0.15	0.20	0.18	0.15	0.19	0.18	0.18	0.16	0.15
$\chi^2_{(10)}$	0.13	0.11	0.12	0.10	0.09	0.12	0.12	0.11	0.11	0.11
Logistic	0.11	0.10	0.08	0.06	0.05	0.09	0.07	0.08	0.09	0.09

Table 8. Empirical power of tests for $\alpha = 0.10$ and $I = 10$

Distribution	W	S2	S12	S25	S26	S13	S61	S28	S40	S43
Binomial	0.55	0.56	0.59	0.60	0.60	0.58	0.60	0.61	0.56	0.54
Uniform	0.17	0.12	0.23	0.28	0.31	0.20	0.25	0.23	0.21	0.21
Cauchy	0.62	0.66	0.61	0.59	0.58	0.62	0.60	0.60	0.60	0.59
Half-normal	0.29	0.27	0.29	0.30	0.29	0.28	0.28	0.26	0.24	0.22
$\chi^2_{(10)}$	0.18	0.18	0.17	0.17	0.16	0.17	0.17	0.17	0.16	0.16
Logistic	0.16	0.16	0.15	0.13	0.12	0.15	0.14	0.14	0.14	0.14

Table 9. Empirical power of tests for $\alpha = 0.03$ and $I = 20$

Distribution	W	S2	S12	S25	S26	S13	S61	S28	S40	S43
$\chi^2_{(1)}$	0.97	0.96	0.97	0.96	0.93	0.98	0.97	0.98	0.97	0.95
$\chi^2_{(2)}$	0.80	0.76	0.76	0.69	0.61	0.77	0.76	0.77	0.73	0.69
$\chi^2_{(4)}$	0.46	0.43	0.42	0.35	0.29	0.43	0.42	0.42	0.37	0.36
$\chi^2_{(10)}$	0.21	0.18	0.19	0.15	0.11	0.19	0.18	0.19	0.17	0.16
Non-cent. $\chi^2_{(16)}$	0.12	0.12	0.10	0.08	0.06	0.11	0.11	0.12	0.11	0.11
Log normal	0.93	0.91	0.91	0.87	0.82	0.91	0.90	0.90	0.87	0.84

Table 9 (Continued)

Distribution	W	S2	S12	S25	S26	S13	S61	S28	S40	S43
Cauchy	0.83	0.85	0.80	0.76	0.72	0.82	0.81	0.82	0.82	0.82
Uniform	0.15	0.07	0.22	0.29	0.31	0.16	0.26	0.21	0.17	0.20
Logistic	0.08	0.09	0.06	0.05	0.03	0.09	0.08	0.10	0.10	0.10
Beta (2, 1)	0.24	0.17	0.27	0.28	0.26	0.23	0.28	0.26	0.24	0.24
LaPlace	0.24	0.27	0.18	0.11	0.07	0.22	0.19	0.23	0.26	0.25
Poisson	1.00	0.99	1.00	1.00	0.98	0.99	0.99	0.98	0.85	0.80
Binomial	0.64	0.61	0.72	0.74	0.70	0.61	0.66	0.60	0.37	0.33
T(5, 2.4)	0.47	0.36	0.52	0.57	0.55	0.48	0.56	0.54	0.52	0.52
T(10, 3.1)	0.83	0.76	0.85	0.84	0.79	0.83	0.86	0.86	0.84	0.80
Half-normal	0.36	0.30	0.35	0.32	0.26	0.34	0.35	0.35	0.34	0.32
Half-Cauchy	0.98	0.98	0.98	0.98	0.97	0.98	0.98	0.98	0.97	0.97
Sum of 3 uniforms	0.03	0.02	0.03	0.04	0.04	0.03	0.03	0.04	0.03	0.03

Table 10. Empirical power of tests for $\alpha = 0.05$ and $I = 20$

Distribution	W	S2	S12	S25	S26	S13	S61	S28	S40	S43
$\chi^2_{(1)}$	0.99	0.98	0.98	0.97	0.96	0.99	0.99	0.99	0.99	0.98
$\chi^2_{(2)}$	0.84	0.83	0.84	0.78	0.72	0.84	0.83	0.84	0.80	0.76
$\chi^2_{(4)}$	0.53	0.53	0.53	0.45	0.40	0.52	0.51	0.52	0.47	0.43
$\chi^2_{(10)}$	0.28	0.29	0.27	0.20	0.18	0.26	0.24	0.25	0.23	0.21
Non-cent. $\chi^2_{(16)}$	0.15	0.16	0.15	0.13	0.11	0.14	0.15	0.15	0.15	0.15
Log normal	0.94	0.94	0.94	0.91	0.88	0.94	0.93	0.93	0.91	0.89
Cauchy	0.85	0.87	0.83	0.79	0.77	0.85	0.83	0.85	0.85	0.84
Uniform	0.22	0.15	0.29	0.36	0.39	0.23	0.32	0.29	0.26	0.27
Logistic	0.11	0.14	0.09	0.06	0.05	0.12	0.10	0.13	0.14	0.13
Beta (2, 1)	0.31	0.27	0.35	0.35	0.34	0.33	0.36	0.36	0.34	0.33
LaPlace	0.28	0.35	0.25	0.16	0.13	0.29	0.24	0.29	0.31	0.29

Table 10 (Continued)

Distribution	W	S2	S12	S25	S26	S13	S61	S28	S40	S43
Poisson	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.99	0.91	0.87
Binomial	0.73	0.74	0.82	0.85	0.83	0.70	0.74	0.67	0.47	0.40
T(5, 2.4)	0.55	0.48	0.63	0.65	0.64	0.58	0.65	0.65	0.64	0.63
T(10, 3.1)	0.88	0.84	0.90	0.88	0.87	0.89	0.91	0.92	0.91	0.89
Half-normal	0.43	0.42	0.44	0.40	0.36	0.44	0.45	0.46	0.43	0.40
Half-Cauchy	0.99	0.98	0.99	0.99	0.98	0.99	0.99	0.99	0.99	0.99
Sum of 3 uniforms	0.05	0.04	0.06	0.07	0.07	0.04	0.06	0.05	0.05	0.05

Table 11. Empirical power of tests for $\alpha = 0.10$ and $I = 20$

Distribution	W	S2	S12	S25	S26	S13	S61	S28	S40	S43
$\chi^2_{(1)}$	0.99	0.99	0.99	0.99	0.98	0.99	0.99	0.99	0.99	0.99
$\chi^2_{(2)}$	0.90	0.89	0.90	0.88	0.86	0.91	0.90	0.91	0.91	0.89
$\chi^2_{(4)}$	0.67	0.65	0.65	0.61	0.56	0.66	0.65	0.64	0.61	0.58
$\chi^2_{(10)}$	0.38	0.35	0.36	0.33	0.31	0.38	0.36	0.36	0.33	0.31
Non-cent. $\chi^2_{(16)}$	0.25	0.25	0.24	0.20	0.20	0.25	0.24	0.24	0.23	0.23
Log normal	0.97	0.96	0.96	0.95	0.94	0.97	0.96	0.97	0.96	0.95
Cauchy	0.88	0.90	0.87	0.84	0.83	0.87	0.86	0.86	0.86	0.86
Uniform	0.37	0.25	0.41	0.49	0.52	0.36	0.44	0.41	0.38	0.41
Logistic	0.18	0.21	0.15	0.11	0.10	0.19	0.16	0.19	0.19	0.19
Beta (2, 1)	0.48	0.41	0.48	0.50	0.49	0.49	0.50	0.50	0.50	0.49
LaPlace	0.39	0.45	0.35	0.29	0.25	0.38	0.34	0.37	0.35	0.34
Poisson	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.97	0.92
Binomial	0.94	0.92	0.99	0.99	1.00	0.86	0.89	0.78	0.63	0.57
T(5, 2.4)	0.73	0.63	0.74	0.76	0.77	0.74	0.78	0.77	0.78	0.78
T(10, 3.1)	0.95	0.92	0.95	0.95	0.93	0.95	0.96	0.96	0.96	0.95
Half-normal	0.59	0.54	0.57	0.55	0.53	0.60	0.60	0.61	0.59	0.57

Table 11 (Continued)

Distribution	W	S2	S12	S25	S26	S13	S61	S28	S40	S43
Half-Cauchy	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.99	1.00	0.99
Sum of 3 uniforms	0.12	0.10	0.12	0.12	0.12	0.11	0.12	0.11	0.10	0.10

Table 12. Empirical power of tests for $\alpha = 0.03$ and $I = 50$

Distribution	W	S2	S12	S25	S26	S13	S61	S28	S40	S43
Binomial	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.35	0.23
Uniform	0.81	0.51	0.77	0.82	0.85	0.70	0.81	0.77	0.78	0.83
Cauchy	0.99	1.00	0.99	0.99	0.99	0.99	0.99	0.99	0.98	0.98
Half-normal	0.92	0.87	0.91	0.88	0.82	0.91	0.92	0.92	0.90	0.84
Logistic	0.11	0.21	0.14	0.07	0.05	0.20	0.16	0.19	0.20	0.18

Table 13. Empirical power of tests for $\alpha = 0.05$ and $I = 50$

Distribution	W	S2	S12	S25	S26	S13	S61	S28	S40	S43
Uniform	0.87	0.65	0.81	0.87	0.90	0.79	0.87	0.85	0.84	0.88
Binomial	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.49	0.38
Cauchy	0.99	1.00	0.99	0.99	0.99	0.99	0.99	0.99	0.99	0.98
Half-normal	0.95	0.92	0.94	0.93	0.90	0.95	0.96	0.96	0.94	0.92
Logistic	0.14	0.28	0.19	0.12	0.09	0.25	0.22	0.24	0.24	0.22

Table 14. Empirical power of tests for $\alpha = 0.10$ and $I = 50$

Distribution	W	S2	S12	S25	S26	S13	S61	S28	S40	S43
Binomial	1.00	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.83	0.68
Uniform	0.95	0.80	0.90	0.93	0.94	0.88	0.93	0.91	0.93	0.95
Cauchy	1.00	1.00	1.00	1.00	1.00	1.00	0.99	0.99	0.99	0.99
Half-normal	0.98	0.96	0.97	0.97	0.96	0.98	0.98	0.98	0.98	0.97
Logistic	0.22	0.34	0.29	0.23	0.19	0.33	0.30	0.31	0.30	0.29

Table 15. Differences between empirical powers of W and S tests for $\alpha = 0.03$ and $I = 10$

Distribution	S2	S12	S25	S26	S13	S61	S28	S40	S43
Binomial	+0.05	+0.13	+0.09	+0.09	+0.10	+0.11	+0.09	+0.01	-0.03
Uniform	-0.02	+0.05	+0.07	+0.09	+0.02	+0.06	+0.04	+0.04	+0.05
Cauchy	+0.01	-0.03	-0.09	-0.13	-0.01	-0.04	-0.03	-0.04	-0.04
Half-normal	-0.01	+0.03	-0.01	-0.03	0	0	0	-0.02	-0.03
$\chi^2_{(10)}$	+0.01	+0.02	-0.01	-0.02	+0.01	0	+0.01	0	0
Logistic	0	-0.02	-0.03	-0.04	-0.01	-0.02	-0.02	-0.02	-0.02

Table 16. Differences between empirical powers of W and S tests for $\alpha = 0.05$ and $I = 10$

Distribution	S2	S12	S25	S26	S13	S61	S28	S40	S43
Binomial	-0.09	+0.02	+0.05	+0.04	+0.02	+0.04	+0.03	0	-0.01
Uniform	-0.04	+0.05	+0.09	+0.10	+0.02	+0.06	+0.04	+0.04	+0.04
Cauchy	+0.01	-0.02	-0.06	-0.09	-0.01	-0.03	-0.03	-0.04	-0.05
Half-normal	-0.03	+0.02	0	-0.03	+0.01	0	0	-0.02	-0.03
$\chi^2_{(10)}$	-0.02	-0.01	-0.03	-0.04	-0.01	-0.01	-0.02	-0.02	-0.02
Logistic	-0.01	-0.03	-0.05	-0.06	-0.02	-0.04	-0.03	-0.02	-0.02

Table 17. Differences between empirical powers of W and S tests for $\alpha = 0.10$ and $I = 10$

Distribution	S2	S12	S25	S26	S13	S61	S28	S40	S43
Binomial	+0.01	+0.04	+0.05	+0.05	+0.03	+0.05	+0.06	+0.01	-0.01
Uniform	-0.05	+0.06	+0.11	+0.14	+0.03	+0.08	+0.06	+0.04	+0.04
Cauchy	+0.04	-0.01	-0.03	-0.04	0	-0.02	-0.02	-0.02	-0.03
Half-normal	-0.02	0	+0.01	0	-0.01	-0.01	-0.03	-0.05	-0.07
$\chi^2_{(10)}$	0	-0.01	-0.01	-0.02	-0.01	-0.01	-0.01	-0.02	-0.02
Logistic	0	-0.01	-0.03	-0.04	-0.01	-0.02	-0.02	-0.02	-0.02

Table 18. Differences between empirical powers of W and S tests for $\alpha = 0.03$ and $I = 20$

Distribution	S2	S12	S25	S26	S13	S61	S28	S40	S43
$\chi^2_{(1)}$	-0.01	0	-0.01	-0.04	+0.01	0	+0.01	0	-0.02
$\chi^2_{(2)}$	-0.04	-0.04	-0.11	-0.19	-0.03	-0.04	-0.03	-0.07	-0.11
$\chi^2_{(4)}$	-0.03	-0.04	-0.11	-0.17	-0.03	-0.04	-0.04	-0.09	-0.10
$\chi^2_{(10)}$	-0.03	-0.02	-0.06	-0.10	-0.02	-0.03	-0.02	-0.04	-0.05
Non-cent. $\chi^2_{(16)}$	0	-0.02	-0.04	-0.06	-0.01	-0.01	0	-0.01	-0.01
Log-normal	-0.02	-0.02	-0.06	-0.11	-0.02	-0.03	-0.03	-0.06	-0.09
Cauchy	+0.02	-0.03	-0.07	-0.11	-0.01	-0.02	-0.01	-0.01	-0.01
Uniform	-0.08	+0.07	+0.14	+0.16	+0.01	+0.11	+0.06	+0.02	+0.05
Logistic	+0.01	-0.02	-0.03	-0.05	+0.01	0	+0.02	+0.02	+0.02
Beta (2, 1)	-0.07	+0.03	+0.04	+0.02	-0.01	+0.04	+0.02	0	0
LaPlace	+0.03	-0.06	-0.13	-0.17	-0.02	-0.05	-0.01	+0.02	+0.01
Poisson	-0.01	0	0	-0.02	-0.01	-0.01	-0.02	-0.15	-0.20
Binomial	-0.03	+0.08	+0.10	+0.06	-0.03	+0.02	-0.04	-0.27	-0.31
T(5, 2.4)	-0.11	+0.05	+0.10	+0.08	+0.01	+0.09	+0.07	+0.05	+0.05
T(10, 3.1)	-0.07	+0.02	+0.01	-0.04	0	+0.03	+0.03	+0.01	-0.03
Half-normal	-0.06	-0.01	-0.04	-0.10	-0.02	-0.01	-0.01	-0.02	-0.04

Table 18 (Continued)

Distribution	S2	S12	S25	S26	S13	S61	S28	S40	S43
Half-Cauchy	0	0	0	-0.01	0	0	0	-0.01	-0.01
Sum of 3 uniforms	-0.01	0	+0.01	+0.01	0	0	+0.01	0	0

Table 19. Differences between empirical powers of W and S tests for $\alpha = 0.05$ and $I = 20$

Distribution	S2	S12	S25	S26	S13	S61	S28	S40	S43
$\chi^2_{(1)}$	-0.01	-0.01	-0.02	-0.03	0	0	0	0	-0.01
$\chi^2_{(2)}$	-0.01	0	-0.06	-0.12	0	-0.01	0	-0.04	-0.08
$\chi^2_{(4)}$	0	0	-0.08	-0.13	-0.01	-0.02	-0.01	-0.06	-0.10
$\chi^2_{(10)}$	+0.01	-0.01	-0.07	-0.10	-0.02	-0.04	-0.03	-0.05	-0.07
Non-cent. $\chi^2_{(16)}$	+0.01	0	-0.02	-0.04	-0.01	0	0	0	0
Log normal		0	-0.03	-0.06	0	-0.01	-0.01	-0.03	-0.05
Cauchy	+0.02	-0.02	-0.06	-0.08	0	-0.02	0	0	-0.01
Uniform	-0.07	+0.07	+0.14	+0.17	+0.01	+0.10	+0.07	+0.04	+0.05
Logistic	+0.03	-0.02	-0.05	-0.06	+0.01	-0.01	+0.02	+0.03	+0.02
Beta (2, 1)	-0.04	+0.04	+0.04	+0.03	+0.02	+0.05	+0.05	+0.03	+0.02
LaPlace	+0.07	-0.03	-0.12	-0.15	+0.01	-0.04	+0.01	+0.03	+0.01
Poisson	0	0	0	0	0	0	-0.01	-0.09	-0.13
Binomial	+0.01	+0.09	+0.12	+0.10	-0.03	+0.01	-0.06	-0.26	-0.33
T(5, 2.4)	-0.07	+0.08	+0.10	+0.09	+0.03	+0.10	+0.10	+0.09	+0.08
T(10, 3.1)	-0.04	+0.02	0	-0.01	+0.01	+0.03	+0.04	+0.03	+0.01
Half-normal	-0.01	+0.01	-0.03	-0.07	+0.01	+0.02	+0.03	0	-0.03
Half-Cauchy	-0.01	0	0	-0.01	0	0	0	0	0
Sum of 3 uniforms	-0.01	+0.01	+0.02	+0.02	-0.01	+0.01	0	0	0

Table 20. Differences between empirical powers of W and S tests for $\alpha = 0.10$ and $I = 20$

Distribution	S2	S12	S25	S26	S13	S61	S28	S40	S43
$\chi^2_{(1)}$	0	0	0	-0.01	0	0	0	0	0
$\chi^2_{(2)}$	-0.01	0	-0.02	-0.04	+0.01	0	+0.01	+0.01	-0.01
$\chi^2_{(4)}$	-0.02	-0.02	-0.06	-0.11	-0.01	-0.02	-0.03	-0.06	-0.09
$\chi^2_{(10)}$	-0.03	-0.02	-0.05	-0.07	0	-0.02	-0.02	-0.05	-0.07
Non-cent. $\chi^2_{(16)}$	0	-0.01	-0.05	-0.05	0	-0.01	-0.01	-0.02	-0.02
Log normal	-0.01	-0.01	-0.02	-0.03	0	-0.01	0	-0.01	-0.02
Cauchy	+0.02	-0.01	-0.04	-0.05	0	-0.02	-0.02	-0.02	-0.02
Uniform	-0.12	+0.04	+0.12	+0.15	-0.01	+0.07	+0.04	+0.01	+0.04
Logistic	+0.03	-0.03	-0.07	-0.08	+0.01	-0.02	+0.01	+0.01	+0.01
Beta (2, 1)	-0.07	0	+0.02	+0.01	+0.01	+0.02	+0.02	+0.02	+0.01
LaPlace	+0.06	-0.04	-0.10	-0.14	-0.01	-0.05	-0.02	-0.04	-0.05
Poisson	0	0	0	0	0	0	0	-0.03	-0.08
Binomial	-0.02	+0.05	+0.05	+0.06	-0.08	-0.05	-0.16	-0.31	-0.37
T(5, 2.4)	-0.10	+0.01	+0.03	+0.04	+0.01	+0.05	+0.04	+0.05	+0.05
T(10, 3.1)	-0.03	0	0	-0.02	0	+0.01	+0.01	+0.01	0
Half-normal	-0.05	-0.02	-0.04	-0.06	+0.01	+0.01	+0.02	0	-0.02
Half-Cauchy	0	0	0	0	0	0	0	+0.01	0
Sum of 3 uniforms	-0.02	0	0	0	-0.01	0	-0.01	-0.02	-0.02

Table 21. Differences between empirical powers of W and S tests for $\alpha = 0.03$ and $I = 50$

Distribution	S2	S12	S25	S26	S13	S61	S28	S40	S43
Binomial	0	0	0	0	0	0	0	-0.65	-0.77
Uniform	-0.30	-0.04	+0.01	+0.04	-0.11	0	-0.04	-0.03	+0.02
Cauchy	0	+0.01	0	0	0	0	0	-0.01	-0.01

Table 21 (Continued)

Distribution	S2	S12	S25	S26	S13	S61	S28	S40	S43
Half-normal	-0.05	-0.01	-0.04	-0.10	-0.01	0	0	-0.02	-0.08
Logistic	+0.10	+0.03	-0.04	-0.06	+0.09	+0.05	+0.08	+0.09	+0.07

Table 22. Differences between empirical powers of W and S tests for $\alpha = 0.05$ and $I = 50$

Distribution	S2	S12	S25	S26	S13	S61	S28	S40	S43
Binomial	0	0	0	0	0	0	0	-0.51	-0.62
Uniform	-0.22	-0.05	0	+0.03	-0.08	0	-0.02	-0.03	+0.01
Cauchy	+0.01	0	0	0	0	0	0	0	-0.01
Half-normal	-0.03	-0.01	-0.02	-0.05	0	+0.01	+0.01	-0.01	-0.03
Logistic	+0.14	+0.05	-0.02	-0.05	+0.11	+0.08	+0.10	+0.10	+0.08

Table 23. Differences between empirical powers of W and S tests for $\alpha = 0.10$ and $I = 50$

Distribution	S2	S12	S25	S26	S13	S61	S28	S40	S43
Binomial	0	0	0	0	0	0	0	-0.17	-0.32
Uniform	-0.15	-0.05	-0.02	-0.01	-0.07	-0.02	-0.04	-0.02	0
Cauchy	0	0	0	0	0	-0.01	-0.01	-0.01	-0.01
Half-normal	-0.02	-0.01	-0.01	-0.02	0	0	0	0	-0.01
Logistic	+0.12	+0.07	+0.01	-0.03	+0.11	+0.08	+0.09	+0.08	+0.07

Table 24. Summary of differences between empirical powers of the W and S tests

	S2	S12	S25	S26	S13	S61	S28	S40	S43
Sum of differences	-1.56	+0.43	-0.62	-1.97	-0.03	+0.57	+0.38	-2.67	-4.01
Maximum positive difference	+0.14	+0.13	+0.14	+0.17	+0.11	+0.11	+0.10	+0.09	+0.08
Maximum negative difference	-0.30	-0.06	-0.13	-0.19	-0.11	-0.05	-0.16	-0.65	-0.77

The values of the S statistics S61, S12, S28, and S13 for $\alpha = 0.03$, 0.05, and 0.10 and sample sizes of 10, 20 and 50 are shown in Table 25. A difference in methods of application should be noted. For the W test, the null hypothesis should be rejected if the computed W is less than the table $W_{\alpha, I}$. For the S tests, the null hypothesis should be rejected if the computed S is greater than the table $S_{\alpha, I}$. The S test procedure is:

1. $H_0: \{y_i\} \sim N(\mu, \sigma)$
2. Select α and sample size, I.
3. Draw sample and arrange observed values in ascending order.
4. Obtain standard deviates of the normal distribution from a table and arrange in ascending order (the continuity correction used in this study was $(2_i - 1)/(2I)$).
5. Fit a linear regression line to the paired $\{X_i, Y_i\}$.
6. Compute $\bar{Y} = \sum_{i=1}^I Y_i / I$ and $\hat{Y}_i = a + bX_i$.
7. Compute S^{**} (S61, S12, S28, or S13).
8. If $S^{**} > S_{\alpha, I}^{**}$, reject H_0 .

Table 25. Values of the S61, S12, S28, and S13 statistics

		S61	S12	S28	S13
I = 10	$\alpha = 0.03$	0.95348E-2 ^a	0.93992E-1	0.50630E-2	0.63203E-2
	$\alpha = 0.05$	0.72391E-2	0.79797E-1	0.39329E-2	0.46956E-2
	$\alpha = 0.10$	0.47298E-2	0.60930E-1	0.24887E-2	0.30499E-2
I = 20	$\alpha = 0.03$	0.44251E-2	0.68788E-1	0.20886E-2	0.25467E-2
	$\alpha = 0.05$	0.34118E-2	0.54936E-1	0.15477E-2	0.18248E-2
	$\alpha = 0.10$	0.21474E-2	0.40838E-1	0.10039E-2	0.10967E-2
I = 50	$\alpha = 0.03$	0.17546E-2	0.37098E-1	0.82431E-3	0.64293E-3
	$\alpha = 0.05$	0.11906E-2	0.31369E-1	0.56386E-3	0.47066E-3
	$\alpha = 0.10$	0.75836E-3	0.22120E-1	0.34697E-3	0.28435E-3

^aThe number before E is to be multiplied by a factor of ten raised to the power of the algebraic number after E.

The results of the final simulation runs indicate:

1. The S61, S12, S28, and S13 tests appear to be the "best" of the S type tests.
2. The empirical powers of the S61, S12, S28, and S13 tests are comparable to the empirical powers of the W tests.

The criterion used in selecting the best of the S tests was the total sum of the differences between the empirical powers of the W and the S tests. An additional criterion was to minimize the maximum negative differences between the powers of the W and S tests.